1 Thank you for the thoughtful feedback and comments; we are delighted to see this positive response. All reviewers agree
2 that we tackle an important problem of interest to the NeurIPS community, and acknowledge our extensive/insightful
3 experimental analysis which shows improvement over state-of-the-art. The reviewers recognize that we present an
4 interesting idea of relying on computational-identifiability [R3] and our proposed adversarial reweighing method is
5 sound [R3] which is composed in a well-written paper [R1,R3,R4]. Thank you! We respond to your questions below:

6 [R1] *When are DRO's benefits apparent ...?* The main benefit of DRO is that it induces distributional robustness.
7 Consequently, DRO is guaranteed to focus on worst-case risk for *any* group in the data exceeding size $\alpha$. In contrast,
8 ARL would only improve the performance for groups that are computationally identifiable groups over (x,y). However,
9 in the fairness setting, this is not necessarily a benefit for DRO, as the DRO authors state in their ICML talk: "it leads to
10 a trade-off between the outlier resistance and fairness, and when the fraction of outliers in training data exceeds size $\alpha$
11 DRO is completely broken." Our experiments on label bias (Fig. 3(b)) shed some light on this. These experiments
12 were performed with DRO parameter $\alpha$ set to 0.2. We observe that while the fraction of incorrect ground truth class
13 labels (i.e., outliers) is less than 0.2, the performance of ARL and DRO is nearly the same. As the fraction of outliers in
14 the training set exceeds 0.2 we observe that DRO's performance drops significantly. We will revise related work, and
15 experiments section to reflect this more clearly.

16 [R1] *The employed reweighting approach is standard, while it offers great performance in the considered problem.*
17 While the reweighted optimization problem is standard, the key contribution of this paper is in the adversarial approach
18 of learning these weights over computable groups in (x,y). Our experiments on "ARL vs Inverse Probability Weighting"
19 (Tbl. 2) illustrate the advantages of ARL over standard re-weighting approaches.

20 [R1] *Is there a unified approach that performs best for all possible scenarios?* Thank you for your thoughtful comment.
21 This along with your other comment on equipping robustness ignited interesting discussions amongst the authors.
22 We think combining the strengths of computationally identifiable regions with distributional robustness would be an
23 interesting research direction to unify the strengths of the two approaches in future research.

24 [R3] *Regarding smoothness of computationally-identifiable regions, and avoiding overfitting to outliers:* Your obser-
25 vations, and understanding is correct. The design and complexity of the adversary model $f_\phi$ plays an important role
26 in controlling the granularity (and smoothness) of computationally-identifiable regions of error. More expressive $f_\phi$
27 lead to finer-grained upweighting but runs the risk of overfitting to outliers. We currently discuss this on line 200-202,
28 but for further clarity we will expand this discussion, and include experimental results on the relationship between the
29 adversary's complexity and the smoothness of the function learnt by the adversary. Further, recall that the adversary
30 example weight function $\lambda_\phi(x_i, y_i) : f_\phi(x_i, y_i) \to \mathbb{R}$ in Eq. 5 is over $f_\phi(x_i, y_i)$. Hence, the distribution of training
31 weights $\lambda_\phi(x_i, y_i)$ is a good indicator of the smoothness of the function $f_\phi$ learnt by the adversary, and a practical
32 heuristic to check if the model is overfitting to training outliers.

33 [R3] *Additional related work:* Thank you for pointing to relevant related work Coston et al. [AIES 2019] and Chen et
34 al. [FAccT 2019]. We will update our related work section. Coston et al. focus on domain adaptation of fairness in
35 settings where the group labels are known for either source or target dataset. Chen et al. investigate the reasons for bias
36 in estimating unfairness in the settings where proxies for group labels are used to evaluate model unfairness.

37 [R3] *Line 184 defines regions as $l(h(x), y) > \epsilon$. But the loss function in Eq. 5 does not adhere to this.* We intended to
38 use this notation from a descriptive perspective as a framing device to explain what we mean by a "significant error".
39 However, this is not necessary for optimization. We will update the model section to clarify this.

40 [R4] *COMPAS protected groups are not computationally-identifiable ... surprising:* While we did not perform a formal
41 study on this, we observed that the adversarial predictive accuracy for race in COMPAS is 0.61 is consistent with prior
42 work [1]. We believe that there are demographic signals in the data, but they are not strong enough to predict groups well.

43 [R4] *How should one understand if computationally-identifiable groups will include the protected groups?* Even
44 evaluating fairness without any group labels is very difficult[2], and as such with *zero* information this would be hard
45 to do. However, we expect in many scenarios practitioners have a variety of domain knowledge they can leverage.
46 For example, if groups labels were not available in training but for a small curated validation set or in an auxiliary
47 dataset, one could evaluate if the computationally identified groups include protected groups, e.g., by analyzing the
48 learnt weights (as in Fig. 4). Practioners might also have knowledge about their domain, e.g., skin tone is present in
49 images and heavily correlated with race. We will add a discussion on limitations of computational identifiability to the
50 paper, and update the broader impact sections to reflect this.

51 [R4] Thank you for helpful additional feedback. We will address all the comments by appropriately revising the text,
52 and reiterating some of the earlier limitations on label noise in broader impact sections for further clarity.

---

[1] See Fig. 4 in "iFair: Learning individually fair data representations for algorithmic decision making. In ICDE 2019. "
[2] See "Assessing algorithmic fairness with unobserved protected class using data combination. In FAccT 2020."