

1 We thank the reviewers for their comments and for acknowledging that we address a relevant problem for the NeurIPS  
2 community [R1,R2,R3], that our experiments show the utility / effectiveness of the proposed method [R1,R2,R3] and  
3 that the setup and theoretical claims seem sound [R1,R2].

4 **R1 Comparison to [11,24] not very fair:** Please note that the claimed missing result “sensitive attribute included”  
5 \*\*\*is in the paper\*\*\* (specifically see the columns in Tabs. 1 & 2 named “Sensitive feature in the functional form of  
6 the model”); See also ll. 233-234. **They do not discuss shifts in  $P(Y|X, S)$ , which will not affect fairness but will  
7 affect accuracy:** The reviewer raises a very interesting point. Note that our method (cf. Eq. (6)) looks for a shared low  
8 complexity representation between the tasks. Shifts in  $P(Y|X, S)$  are permitted, provided there is a common predictive  
9 representation for the task outputs. Searching for such representation is indeed one of the main ideas behind transfer  
10 learning (see e.g. Caruana 1997; Argyriou et al. 2008, etc.) and it is aided by the first term in Eq. (6), measuring the  
11 average empirical risk on the training tasks. However, differently from the above works, our method also crucially  
12 involves the second term, which encourages representations that approximately satisfies DP on average over the training  
13 tasks. Our fairness violation bound (Thm 1) involves only the marginal distributions (of each sensitive group within the  
14 task) since we care to measure this at the representation level, but see ll. 97-103 and Lemma 3 for how this affects DP at  
15 the output level. In this sense our method learn a shared and transferable representation, one based on which accurate  
16 and fair models can be learned on tasks sampled from the environment  $\rho$ . **Would we expect accuracy to degrade at  
17 all?** Yes, accuracy decreases due to the fairness constraint. This property is common to any algorithmic fairness method.  
18 **Other suggestions:** These will be addressed in the revision. In particular we’ll: make our code publicly available upon  
19 acceptance, add at l. 140 that  $\sigma$  is applied component-wise, say at l. 159 that the function “Gap” is quantified by the  
20 r.h.s. of Ineq. (10), clarify at l. 163 that “approximately” means that the marginal distributions of the two sensitive  
21 groups within the task are closed according to some suitable measure (e.g. see that at l. 242), improve the “M [24]”  
22 notation (it stays for “method [24]” but we could just write “[24]”). Finally concerning **fairness/accuracy tradeoff in  
23 the generalization setting:** standard bounds from the learning-to-learn literature (e.g. [3]) can be readily used to bound  
24 the risk (or accuracy) of the representation  $h$  found by our method on future tasks by the minimal multitask empirical  
25 risk (over the specifications  $g_t$ ). Now since our method in Eq. (6) minimizes a tradeoff between the multitask empirical  
26 risk and fairness violation, the larger  $\gamma$  the larger the former term, so risk on future tasks will reflect this tradeoff too.

27 **R2 Submission focuses entirely on dem. parity:** We agree DP is not the ultimate fairness notion. Still it is frequently  
28 studied in the literature and our study is a valuable starting point for fair representation learning within the multitask  
29 setting. **May techniques be deployed with good conscious?** Yes, our theoretical and experimental results give an  
30 indication that the method could be valuable in practice and safely deployed. Of course more experiments would be  
31 needed to assess its robustness on on real-world problems. **Fictitious multitask setting:** it is true the paper is more on  
32 theory and methodology which is not always close to practice, but imagine the following real-life scenario: each task is  
33 associated with an hospital in country X and the task is to predict whether a patient who visits emergency should be  
34 hospitalized. The environment (meta-distribution)  $\rho$  may be the uniform distribution (or weight larger hospitals more).  
35 The sensitive attribute may be race and other non-sensitive variables may measure cough frequency and body temperature.  
36 Our main result, Thm 1 (in conjunction with Cor 2) then says that if we use our method to learn a predictive representation  
37 and observe it to be fair according to DP on the random training task datasets, then it will also be fair according to DP  
38 on average on all possible hospitals at the population level (i.e. on average over random patients visiting the hospital),  
39 which is a very appealing property. **Description of  $\rho$ :** Yes, we’ll mentioning also at l. 104. **How is the multitask  
40 setting realized in the empirical experiments?** Please see ll. 227-34: we test either on different data for the same  
41 tasks using during training or on a new task in leave-one-task out setting. **How is the Sinkhorn divergence in eq. (5)  
42 estimated from finite samples?** Consider  $\hat{P}$  and  $\hat{Q}$  the same as in ll. 116-117. Denote by  $\mathbf{p} = (p_1, \dots, p_n) := 1/n\mathbf{1}_n$   
43 and  $\mathbf{q} = (q_1, \dots, q_m) := 1/m\mathbf{1}_m$ , with  $\mathbf{1}_k$  the vector with  $k$  entries equal to one ( $\mathbf{p}$  and  $\mathbf{q}$  denote the weights of  
44 the empirical distributions  $\hat{P}$  and  $\hat{Q}$ ). Then,  $\text{OT}_\varepsilon(\hat{P}, \hat{Q}) = \min_{T \in \Pi(\mathbf{p}, \mathbf{q})} \langle T, C \rangle + \varepsilon \sum_{i,j=1}^{n,m} \log(T_{ij}/p_i q_j) T_{ij}$ , where  
45  $\Pi(\mathbf{p}, \mathbf{q}) := \{T \in \mathbb{R}_+^{n \times m} \mid T\mathbf{1}_m = \mathbf{p}, T\mathbf{1}_n = \mathbf{q}\}$  and  $C_{ij} = \|x_i - z_j\|^2$ ; see ref. [30] for more explanations. **Restriction  
46 to 1-hidden layer nets:** Our method and Thm 1 apply to general classes of representation functions  $h$  of suitably  
47 bounded complexity. For simplicity we illustrated them on 1-hidden layer networks, both theoretically (Thm 1 + Cor 2)  
48 and empirically. However, bounds in [Chain Rule for the Expected Suprema of Gaussian Processes, ALT 2014] could  
49 be used in place of Cor 2 for multi-layer representations. **Extending Thm 1 to Sinkhorn divergence (l. 188):** There  
50 are two main obstacles at this stage: first, the term  $A_h$  at l. 425 would not be zero, because the estimators is biased.  
51 Second, the Lipschitz behaviour needed to factor out the constant (see l. 436) is not clear in the case of Sinkhorn.

52 **R3 1. Main contribution:** Please see l. 46. To the best of our knowledge this is the first paper using multitask learning  
53 for fair representation. **2. Significance of Thm 1:** This result together with a bound on the Rademacher average of the  
54 representation class (e.g. Cor. 2 in the case of linear representations) gives a justification for our method – see also the  
55 reply to R2 concerning the point *fictitious multitask setting*. **3. Contribution insufficient:** We disagree: MMD and  
56 SNK have been used only recently for algorithmic fairness. The proposed method is novel, empirically competitive and  
57 theoretically grounded (see point 2 above). **4.5. Formula mistakes and missing definitions:** We’ll carefully check  
58 our formulas. Thanks for Eq. (2), but note  $d(\cdot)$  is already defined at l. 109 of the paper.