

1 **Author Response for Submission 8969: Differential Privacy Has Disparate Impact on Model Accuracy**

2 **Related work.** Thanks for the pointers to recently released papers, we will acknowledge them. Yeom et al. [5] show
3 that models with poor generalization are more vulnerable to inference attacks. They measure how DP bounds the
4 leakage of information about training data, not its impact on model accuracy. They do not at all study (1) the accuracy
5 of models on subgroups, nor (2) how accuracy on subgroups changes as a result of applying DP-SGD. (2) is our main
6 result, which is completely independent and orthogonal to the analysis in Yeom et al.

7 We did not have room to discuss [3, 2] but the brief summary is they provide evidence that DP may be combined with
8 fairness, but do not give algorithms that could be used to train practical DP neural networks.

9 **Rényi differential privacy.** We use Rényi DP only to estimate privacy loss. This does not change the DPSGD algorithm
10 of Abadi et al. but rather provides tighter bounds on privacy loss [4], allowing to reduce the amount of added noise. The
11 TF Privacy tool enables estimation of epsilon given the input parameters (dataset size, number of epochs, batch size,
12 noise, delta) before starting the training, thus this computation is not part of Algorithm 1. We ensure that our training
13 uses the same hyperparameters as used to estimate epsilon.

14 **General statements about DP.** We will clarify in the abstract and intro that our results apply to DPSGD, a popular
15 way to train DP neural networks, and not necessarily to DP as a general concept.

16 **Experiment details.** Thanks for the comments about improving presentation (captions and trend lines). We used the
17 UTK dataset as an additional source of darker-skinned faces because in the DiF dataset, some individuals with lighter
18 skin were labeled as dark-skinned. We set the ratio between lighter- and darker-skinned individuals to measure the
19 effect of DPSGD on underrepresented classes, not to reflect the demographic balance of any country or real-world
20 dataset.

21 **Size of the groups and complex classes.** More items per class is indeed usually helpful. That said, our federated
22 learning study shows that participants with simpler vocabularies get better accuracy with DPSGD, whereas participants
23 with complex vocabularies contribute less to the model (Figure 3b). This is an example of how DPSGD negatively
24 affects more complex data.

25 **Impact of clipping.** Clipping alone is responsible for slowing down the learning, similar to decreasing the learning
26 rate. Without adding noise, both well- and under-represented classes converge to the same accuracy but much slower.
27 Noise, however, prevents the model from converging to the same norm. We find this presentation to be more intuitive
28 and perhaps a good starting point for future research on combining differential privacy with fairness.

29 **Adversarial training.** Adversarial training for fairness [1] overweights the loss for underrepresented groups. Sensitivity
30 bounds imposed by DPSGD, DPGAN, and similar approaches hold only for specific loss functions and sampling
31 strategies; if combined directly with adversarial training, the resulting models will not be DP. It is an open problem how
32 to combine DP with censoring techniques such as adversarial training.

33 **Training models with the same epsilon.** Modifying the hyperparameters directly involved in estimating epsilon results
34 in a big variance of results. Using TF Privacy, we observed that among all hyperparameters, the noise multiplier z
35 has the highest impact on epsilon. Changing hyperparameters that do not affect privacy loss, such as the learning rate,
36 model architecture, or optimizer, impacts the accuracy but does not affect fairness, thus we omitted these analyses due
37 to lack of space.

38 **Fairness measure.** Equalized odds gives us the most direct way to measure the impact of DPSGD on a popular
39 fairness measure. Equal opportunity requires equality on the "advantaged" outcome, but in the multi-label tasks in our
40 experiments it is not always clear what outcome should be considered advantaged. Accuracy on each subgroup, on the
41 other hand, is straightforward to measure.

42 REFERENCES

- 43 [1] A. Beutel, J. Chen, Z. Zhao, and E. H. Chi. Data decisions and theoretical implications when adversarially learning
44 fair representations. *arXiv:1707.00075*, 2017.
- 45 [2] R. Cummings, V. Gupta, D. Kimpara, and J. Morgenstern. On the compatibility of privacy
46 and fairness. [http://wpw.gatech.edu/rachel-cummings/wp-content/uploads/sites/679/2019/03/
47 FairPrivate.pdf](http://wpw.gatech.edu/rachel-cummings/wp-content/uploads/sites/679/2019/03/FairPrivate.pdf), 2019.
- 48 [3] M. Jagielski, M. Kearns, J. Mao, A. Oprea, A. Roth, S. Sharifi-Malvajerdi, and J. Ullman. Differentially private fair
49 learning. *arXiv:1812.02696*, 2018.
- 50 [4] I. Mironov. Rényi differential privacy. In *CSF*, 2017.
- 51 [5] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha. Privacy risk in machine learning: Analyzing the connection to
52 overfitting. In *CSF*, 2018.