# Appendix for "Poisoning for Debiasing: Fair Recognition via Eliminating Bias Uncovered in Data Poisoning"

**Table 2: The Model Bias (in %, Equalodds, ↓), Avg. Group Accuracy (in %, ↑), and Worst Group Accuracy (in %, ↑) of models trained on CelebA. Here *bn, a, bl* and *bu* respectively denote *bignose, attractive, blonde, bags-under-eyes*. The best results with unknown biases are highlighted in bold. * indicates that the method knows the bias label of training samples.**

| Method | T=bn | | | T=a | | | T=bl | | | T=bu | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | ModelBias↓ | AvgACC↑ | WorstACC↑ | ModelBias↓ | AvgACC↑ | WorstACC↑ | ModelBias↓ | AvgACC↑ | WorstACC↑ | ModelBias↓ | AvgACC↑ | WorstACC↑ |
| Vanilla | $31.40_{\pm0.6}$ | $71.01_{\pm0.3}$ | $43.36_{\pm0.7}$ | $23.51_{\pm0.3}$ | $74.13_{\pm0.2}$ | $62.34_{\pm0.4}$ | $32.05_{\pm0.3}$ | $82.83_{\pm0.2}$ | $58.32_{\pm0.5}$ | $17.70_{\pm0.4}$ | $71.51_{\pm0.3}$ | $43.71_{\pm0.6}$ |
| Focal | $23.29_{\pm0.5}$ | $70.99_{\pm0.2}$ | $47.61_{\pm0.4}$ | $24.25_{\pm0.4}$ | $76.56_{\pm0.2}$ | $61.36_{\pm0.5}$ | $30.05_{\pm0.3}$ | $83.78_{\pm0.2}$ | $64.35_{\pm0.4}$ | $16.31_{\pm0.3}$ | $70.87_{\pm0.3}$ | $42.47_{\pm0.6}$ |
| LfF | $16.73_{\pm0.4}$ | $68.42_{\pm0.5}$ | $53.12_{\pm0.6}$ | $17.53_{\pm0.7}$ | $76.57_{\pm0.4}$ | $67.36_{\pm0.9}$ | $29.75_{\pm0.6}$ | $75.32_{\pm0.6}$ | $49.10_{\pm1.1}$ | $18.73_{\pm0.4}$ | $70.53_{\pm0.5}$ | $43.51_{\pm0.8}$ |
| JTT | $14.29_{\pm0.5}$ | $72.31_{\pm0.3}$ | $55.09_{\pm0.7}$ | $15.06_{\pm0.5}$ | $77.61_{\pm0.4}$ | $65.33_{\pm0.8}$ | $13.07_{\pm0.5}$ | $85.02_{\pm0.2}$ | $75.53_{\pm0.5}$ | $15.34_{\pm0.4}$ | $70.22_{\pm0.2}$ | $54.01_{\pm0.3}$ |
| DebiAN | $29.03_{\pm0.6}$ | $69.41_{\pm0.3}$ | $39.75_{\pm0.8}$ | $22.39_{\pm0.5}$ | $76.63_{\pm0.3}$ | $59.22_{\pm0.9}$ | $29.35_{\pm0.3}$ | $77.29_{\pm0.3}$ | $63.81_{\pm0.7}$ | $19.38_{\pm0.4}$ | $70.65_{\pm0.3}$ | $44.95_{\pm0.5}$ |
| Echoes | $19.95_{\pm0.6}$ | $66.19_{\pm0.4}$ | $42.73_{\pm0.9}$ | $27.52_{\pm0.4}$ | $72.57_{\pm0.3}$ | $65.11_{\pm0.8}$ | $18.27_{\pm0.5}$ | $76.53_{\pm0.2}$ | $63.52_{\pm0.7}$ | $16.52_{\pm0.7}$ | $70.52_{\pm0.4}$ | $51.23_{\pm0.8}$ |
| BE | $15.57_{\pm0.9}$ | $69.57_{\pm0.5}$ | $56.54_{\pm1.3}$ | $16.21_{\pm0.6}$ | $76.74_{\pm0.4}$ | $67.14_{\pm0.7}$ | $15.84_{\pm0.5}$ | $80.58_{\pm0.4}$ | $69.21_{\pm0.7}$ | $14.32_{\pm0.3}$ | $71.65_{\pm0.6}$ | $59.14_{\pm0.8}$ |
| Poisoner | $\mathbf{6.56_{\pm0.5}}$ | $\mathbf{74.49_{\pm0.2}}$ | $\mathbf{69.61_{\pm0.5}}$ | $\mathbf{3.57_{\pm0.3}}$ | $\mathbf{79.98_{\pm0.2}}$ | $\mathbf{77.04_{\pm0.4}}$ | $\mathbf{8.05_{\pm0.4}}$ | $\mathbf{91.37_{\pm0.3}}$ | $\mathbf{85.78_{\pm0.4}}$ | $\mathbf{7.61_{\pm0.3}}$ | $\mathbf{76.58_{\pm0.2}}$ | $\mathbf{68.29_{\pm0.3}}$ |
| GroupDRO* | $5.54_{\pm0.3}$ | $74.97_{\pm0.2}$ | $66.28_{\pm0.4}$ | $4.02_{\pm0.2}$ | $79.72_{\pm0.2}$ | $76.83_{\pm0.3}$ | $7.61_{\pm0.2}$ | $92.51_{\pm0.1}$ | $81.97_{\pm0.2}$ | $7.96_{\pm0.3}$ | $77.81_{\pm0.1}$ | $67.00_{\pm0.4}$ |

**Table 1: Performance of Poisoner* on the CelebA (T=*bn*) dataset.**

| Method | Avg Acc (↑) | Worst Acc (↑) | Bias (↓) | Bias Accuracy (↑) |
| --- | --- | --- | --- | --- |
| *Vanilla* | $71.01_{\pm0.3}$ | $43.36_{\pm0.7}$ | $31.40_{\pm0.6}$ | - |
| *Poisoner* | $74.49_{\pm0.2}$ | $69.61_{\pm0.5}$ | $6.56_{\pm0.5}$ | $92.26_{\pm0.5}$ |
| *Poisoner** | $74.61_{\pm0.1}$ | $70.53_{\pm0.4}$ | $6.39_{\pm0.6}$ | $93.34_{\pm0.6}$ |

## A  CORRECTIONS TO THE MAIN PAPER

We apologize for the use of two different forms to represent the same concept in the main paper. Specifically, the symbol $D_{con}$ in Lines 378 and 387 should be corrected to $d_{con}$. We apologize for any reading difficulties this may have caused.

## B  ANOTHER METHOD FOR IDENTITY BIAS-CONFLICTING SAMPLES

In the main paper:

> For each sample $x_i$ in the current batch, we compare $x_i$ with the other data in the mini-batch $\mathcal{B}$ as follows:
>
> $$d_{con}(x_i; \theta) = -\frac{1}{|J(i)|} \sum_{j \in J(i)} log \frac{\exp(z_i \cdot z_j)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a)} \quad (1)$$
>
> The larger the $d_{con}(x_i; \theta)$, the more likely it is that the sample $x_i$ is a bias-conflicting sample of the current target class.

The above approach identifies bias-conflicting samples in the model's training on mini-batch $\mathcal{B}$, offering a resource-efficient advantage.

However, this method requires at least two samples with the same label in $\mathcal{B}$, which might not be feasible in cases with a large number of labels (e.g., 1000 classes). To address this situation, we propose using comparison data from outside the mini-batch $\mathcal{B}$.

Specifically, we sampled a larger batch $\mathcal{B}'$ to support comparisons with samples in $\mathcal{B}$. Thus, the calculation of $d_{con}(x_i; \theta)$ for samples in $\mathcal{B}$ is based on both $\mathcal{B}$ and $\mathcal{B}'$.

We implemented our new approach, called Poisoner*, on CelebA (T=*bn*) and present comparative results in Table 1, demonstrating that Poisoner* offers debiasing performance consistent with Poisoner.

**Table 3: The debiasing performance on four benchmark datasets with general bias. * indicates that the method knows the bias label of training samples. - denotes that the test set is not applicable for evaluating model bias.**

| Method | WaterBirds | | | Dogs & Cats | | | C-MNIST[1] | | | C-MNIST[2] | | | ImageNet-B | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ModelBias ↓ | AvgACC ↑ | WorstACC ↑ | ModelBias ↓ | AvgACC ↑ | WorstACC ↑ | ModelBias ↓ | AvgACC ↑ | WorstACC ↑ | ModelBias ↓ | AvgACC ↑ | WorstACC ↑ | ModelBias ↓ | AvgACC ↑ | WorstACC ↑ |
| Vanilla | $35.71_{\pm0.5}$ | $78.17_{\pm0.4}$ | $46.61_{\pm0.7}$ | - | $50.35_{\pm1.3}$ | $47.76_{\pm1.2}$ | $43.27_{\pm2.7}$ | $56.28_{\pm2.8}$ | $10.30_{\pm3.3}$ | $32.73_{\pm2.8}$ | $67.09_{\pm1.3}$ | $7.02_{\pm1.7}$ | $67.09_{\pm0.7}$ | $61.37_{\pm0.3}$ | $23.60_{\pm0.9}$ |
| Focal | $30.12_{\pm0.5}$ | $77.99_{\pm0.3}$ | $56.71_{\pm0.6}$ | - | $68.25_{\pm1.0}$ | $65.85_{\pm1.2}$ | $43.82_{\pm1.7}$ | $56.91_{\pm1.3}$ | $3.05_{\pm1.2}$ | $29.61_{\pm0.2}$ | $70.10_{\pm0.2}$ | $9.01_{\pm0.4}$ | $66.99_{\pm0.7}$ | $61.16_{\pm0.4}$ | $19.30_{\pm0.8}$ |
| LfF | $13.57_{\pm0.5}$ | $79.64_{\pm0.3}$ | $63.56_{\pm0.6}$ | - | $72.91_{\pm0.5}$ | $50.10_{\pm0.6}$ | $12.99_{\pm0.8}$ | $75.26_{\pm0.7}$ | $28.41_{\pm0.9}$ | $11.29_{\pm1.0}$ | $83.92_{\pm0.5}$ | $59.57_{\pm1.2}$ | $37.97_{\pm0.7}$ | $64.21_{\pm0.5}$ | $28.20_{\pm0.8}$ |
| JTT | $12.06_{\pm0.8}$ | $75.61_{\pm0.4}$ | $58.19_{\pm0.9}$ | - | $73.95_{\pm0.5}$ | $67.45_{\pm0.6}$ | $13.07_{\pm0.6}$ | $74.05_{\pm0.5}$ | $30.03_{\pm0.8}$ | $13.38_{\pm0.4}$ | $78.34_{\pm0.3}$ | $56.10_{\pm0.5}$ | $35.24_{\pm0.7}$ | $65.38_{\pm0.6}$ | $30.54_{\pm0.8}$ |
| DebiAN | $12.36_{\pm0.6}$ | $77.72_{\pm0.5}$ | $59.22_{\pm0.7}$ | - | $71.24_{\pm0.5}$ | $67.68_{\pm0.7}$ | $17.08_{\pm0.4}$ | $70.48_{\pm0.4}$ | $36.18_{\pm0.5}$ | $15.25_{\pm0.6}$ | $79.51_{\pm0.5}$ | $57.61_{\pm0.8}$ | $44.58_{\pm0.6}$ | $62.40_{\pm0.5}$ | $29.32_{\pm0.6}$ |
| Echoes | $14.52_{\pm0.5}$ | $78.79_{\pm0.2}$ | $62.73_{\pm0.5}$ | - | $84.56_{\pm0.4}$ | $82.17_{\pm0.8}$ | $16.27_{\pm0.4}$ | $79.18_{\pm0.3}$ | $36.28_{\pm0.5}$ | $13.48_{\pm0.6}$ | $78.42_{\pm0.3}$ | $57.23_{\pm0.7}$ | $38.64_{\pm0.8}$ | $62.54_{\pm0.5}$ | $27.91_{\pm0.8}$ |
| BE | $13.21_{\pm0.4}$ | $76.74_{\pm0.3}$ | $67.14_{\pm0.5}$ | - | $\mathbf{85.59}_{\pm0.3}$ | $79.21_{\pm0.6}$ | $14.98_{\pm0.6}$ | $81.39_{\pm0.4}$ | $\mathbf{39.41}_{\pm0.8}$ | $11.18_{\pm0.8}$ | $85.66_{\pm0.4}$ | $59.14_{\pm0.7}$ | $42.61_{\pm0.4}$ | $64.72_{\pm0.4}$ | $28.96_{\pm0.5}$ |
| Poisoner | $\mathbf{3.57}_{\pm0.3}$ | $\mathbf{84.26}_{\pm0.1}$ | $\mathbf{78.01}_{\pm0.4}$ | - | $84.81_{\pm0.3}$ | $\mathbf{83.02}_{\pm0.5}$ | $\mathbf{12.98}_{\pm0.5}$ | $\mathbf{82.57}_{\pm0.2}$ | $37.40_{\pm0.4}$ | $\mathbf{8.77}_{\pm0.3}$ | $\mathbf{87.51}_{\pm0.1}$ | $\mathbf{64.75}_{\pm0.2}$ | $\mathbf{31.51}_{\pm0.4}$ | $\mathbf{66.28}_{\pm0.3}$ | $\mathbf{32.43}_{\pm0.7}$ |
| GroupDRO* | $5.23_{\pm0.2}$ | $86.72_{\pm0.1}$ | $79.83_{\pm0.2}$ | - | $81.53_{\pm0.4}$ | $68.47_{\pm0.5}$ | $16.51_{\pm0.5}$ | $83.13_{\pm0.4}$ | $29.61_{\pm0.7}$ | $10.58_{\pm0.4}$ | $85.03_{\pm0.2}$ | $42.70_{\pm0.4}$ | $29.22_{\pm0.5}$ | $67.76_{\pm0.3}$ | $44.32_{\pm0.7}$ |

## C STANDARD DEVIATION OF MAIN PAPER

Due to space limitations in the main paper, we did not report the standard deviation. In this appendix, we supplement the experimental results of the main paper by providing the standard deviation in Table 2 and Table 3.