# FAIRDD: FAIR DATASET DISTILLATION VIA SYNCHRONIZED MATCHING

Anonymous authors

Paper under double-blind review

### ABSTRACT

Condensing large datasets into smaller synthetic counterparts has demonstrated its promise for image classification. However, previous research has overlooked a crucial concern in image recognition: ensuring that models trained on condensed datasets are unbiased towards protected attributes (PA), such as gender and race. Our investigation reveals that dataset distillation (DD) fails to alleviate the unfairness towards minority groups within original datasets. Moreover, this bias typically worsens in the condensed datasets due to their smaller size. To bridge the research gap, we propose a novel fair dataset distillation (FDD) framework, namely FairDD, which can be seamlessly applied to diverse matching-based DD approaches, requiring no modifications to their original architectures. The key innovation of FairDD lies in synchronously matching synthetic datasets to PA-wise groups of original datasets, rather than indiscriminate alignment to the whole distributions in vanilla DDs, dominated by majority groups. This synchronized matching allows synthetic datasets to avoid collapsing into majority groups and bootstrap their balanced generation to all PA groups. Consequently, FairDD could effectively regularize vanilla DDs to favor biased generation toward minority groups while maintaining the accuracy of target attributes. Theoretical analyses and extensive experimental evaluations demonstrate that FairDD significantly improves fairness compared to vanilla DD methods, without sacrificing classification accuracy. Its consistent superiority across diverse DDs, spanning Distribution and Gradient Matching, establishes it as a versatile FDD approach.

033

003

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

### 1 INTRODUCTION

034 Deep learning has witnessed remarkable success in computer vision, particularly with recent breakthroughs in vision models (Oquab et al., 2023; Kirillov et al., 2023; Radford et al., 2021; Li et al., 2022; Zhou et al., 2023). Their vision backbones, such as ResNet (He et al., 2015) and ViT (Dosovitskiy et al., 2020), are data-hungry models that require extensive amounts of data for optimization. 037 Dataset Distillation (DD) (Wang et al., 2018; Zhao & Bilen, 2021; 2023; Cazenavette et al., 2022; Wang et al., 2022; Lee et al., 2022b; Cui et al., 2022; Loo et al., 2023; Guo et al., 2023; He & Zhou, 2024; Zhao et al., 2024) provides a promising solution to alleviate this data requirement by 040 condensing the original large dataset into more informative and smaller counterparts (Mehrabi et al., 041 2021; Chung et al., 2023). Despite its appeal, existing DD researches focus on ensuring that models 042 trained on condensed datasets perform comparable accuracy to those trained on the original dataset in 043 terms of target attributes (TA) (Cui et al., 2024; Lu et al., 2024; Vahidian et al., 2024). However, they 044 have overlooked enabling the fairness of trained models with respect to protected attributes (PA).

Unfairness typically arises from imbalanced sample distributions among PA in the empirical training datasets. When the original datasets suffer from the PA imbalance, the corresponding datasets condensed by vanilla DDs inherit and amplify this bias in Fig. 1(a). Since vanilla DDs tend to cover TA distribution for image classification, and as a result, it naturally leads to more synthetic samples located in majority groups compared to minority groups w.r.t. PA, as shown in Fig. 1(b). In this case, these condensed datasets retain the imbalance between protected attributes, thereby rendering the model trained on them unfair. Moreover, the reduced size of the condensed datasets typically amplifies the bias present in the original datasets, especially when there is a significant gap in size between the original and condensed datasets, such as image per class (IPC) 1000 vs. 10. Therefore, it is worthwhile to broaden the scope of DD to encompass both TA accuracy and PA fairness. Recent

072

098

099



073 works attempt to address the class/TA-level long-tailed phenomenon Zhao et al. (2024) and spurious 074 correlations (Cui et al., 2024) to improve the classification performance (Vahidian et al., 2024; Wang et al., 2024), but the exploration on Fair Dataset Distillation (FDD) is still blank. 075

of minority groups and mitigates color imbalance in the condensed dataset via DM+FairDD.

076 To the best of our knowledge, we are the first to propose a novel FDD framework, FairDD, which 077 enables PA fairness in the model trained on the condensed datasets, regardless of PA imbalance in 078 the original data. However, simultaneously maintaining TA accuracy and improving PA fairness is 079 challenging, as the algorithm must properly balance emphasis across all groups, i.e., majority groups 080 primarily for TA distributional coverage and minority groups primarily for PA bias mitigation.

081 FairDD tackles this challenge by (1) partitioning the empirical training distribution into different 082 groups according to PA and decomposing the single alignment target of vanilla DDs into PA-wise 083 sub-targets. (2) synchronously matching synthetic samples to these PA groups, which equally 084 bootstraps synthetic datasets to each PA group without involving the specific group size. In doing 085 so, we reformulated the optimization objectives of vanilla DDs into fairDD-style versions. This allows FairDD to mitigate the effect of the imbalanced PA to the generation of S and to prevent S087 from collapsing into the majority group. Meanwhile, FairDD could also achieve the comprehensive coverage of the entire distribution for TA accuracy as shown in Fig. 1(c). In addition, FairDD requires 088 no modification to existing DDs or additional modules. We provide a theoretical guarantee that FairDD could improve PA fairness while maintaining TA accuracy. 090

091 Extensive experiments demonstrate that our proposed framework effectively mitigates the unfairness 092 in datasets of highly diverse bias. FairDD substantially improves model fairness on condensed datasets compared to various vanilla DDs. FairDD demonstrates its versatility across diverse DD paradigms, including Distribution Matching (DM) and Gradient Matching (GM). Note that we do 094 not apply FairDD to Trajectory Matching (TM) as doing so would require extra model trajectories 095 trained on minority groups, in which the corresponding models would suffer from overfitting due to 096 their limited samples. Our paper makes the following main contributions: 097

- To the best of our knowledge, our research is the first attempt to incorporate fairness issues into DD explicitly. We reveal that vanilla DDs fail to mitigate the bias in original datasets and may exacerbate it due to the limited synthetic samples, leading to severe PA bias in the model trained by the resulting condensed dataset.
- We introduce a novel FDD framework called FairDD, which proposes synchronized match-102 ing to align synthetic samples to all PA groups partitioned from the original data distribution. 103 This allows the generated synthetic samples to be agnostic to PA imbalance of original datasets and maintain the overall distributional coverage of TA. 105
- FairDD requires no alternations to the original architectures of vanilla DDs, in which it only needs to revise the vanilla optimization objectives as a group-level formulation. We also 107 provide theoretical analyses to guarantee the fairness and accuracy of synthetic samples.

109 110

111

112

• Extensive empirical experiments demonstrate that FairDD mitigates the unfairness of vanilla DDs by a large margin. The consistent superiority across diverse DDs, including DM and GM, illustrates that FairDD is a generalist for FDD.

#### 2 PRELIMINARIES

**Dataset Distillation.** Given a vast dataset  $\mathcal{T} = \{(x_i, y_i)\}_{i=1}^N$ , the goal of DD is to condense the original dataset  $\mathcal{T}$  into a smaller dataset  $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^M$  via distillation algorithm Alg with a nerual network, parameterized by  $\theta$ . Randomly initialized classification network  $g_{\psi}$  should maintain the same empirical risk whether it is trained on  $\mathcal{S}$  or  $\mathcal{T}$ .

127

 $S^* = \operatorname*{arg\,min}_{S} \operatorname{Alg}(S, \mathcal{T}, \theta), \tag{1}$ 

$$\mathbb{E}_{\psi \sim \Psi}[\ell(g_{\psi}; \mathcal{S})] \simeq \mathbb{E}_{\psi \sim \Psi}[\ell(g_{\psi}; \mathcal{T})], \tag{2}$$

where Ψ and ℓ(·) represent the parameter space and loss function, respectively. The pioneering
work (Wang et al., 2018) formulates Alg as a bi-level optimization problem. However, such an
optimization process is time-consuming and unsTable Recent works circumvent it and propose
surrogate matching objectives to achieve comparable and even better performance. This research line
is collectively referred to as the DMF, and our paper primarily studies one-stage GM (Zhao et al.,
2020; Zhao & Bilen, 2021) and DM (Zhao & Bilen, 2023; Wang et al., 2022; Zhao & Bilen, 2022).
We leave two-stage trajectory-matching (TM) for future exploration.

**Visual Fairness** Visual fairness is an important field to mitigate discrimination against minority groups. Group fairness requires no statistical disparity to different groups in terms of PA, such as race and gender. This means that an ideal fair model should make independent predictions between TA and PA. One of the common fairness criteria is equalized odds (EO), which computes the prediction accuracy of PA conditioned on TA, to evaluate the level of conditional independence between PA and TA. We use two types of difference of equalized odds DEO<sub>M</sub> and DEO<sub>A</sub> from the worst and averaged levels. Formally, given the PA set  $\mathcal{A} = \{a_1, a_2, ..., a_p\}$ , DEO<sub>M</sub> and DEO<sub>A</sub> (Jung et al., 2021) can be formulated mathematically as follows:

$$\begin{aligned} & \mathsf{DEO}_{\mathsf{M}} = \max_{y \in \mathcal{Y}} \max_{a_i, a_j \in \mathcal{A} \& a_i \neq a_j} \left| P(\hat{Y} = y | Y = y, A = a_i) - P(\hat{Y} = y | Y = y, A = a_j) \right|, \\ & \mathsf{DEO}_{\mathsf{A}} = \max_{y \in \mathcal{Y}} \max_{a_i, a_j \in \mathcal{A} \& a_i \neq a_j} \left| P(\hat{Y} = y | Y = y, A = a_i) - P(\hat{Y} = y | Y = y, A = a_j) \right|. \end{aligned}$$

136 137 138

139

147 148

152

135

## 3 A CLOSE LOOK AT DATASET DISTILLATION FROM FAIRNESS

A unified perspective for Data Match Framework. The essence of the DMF lies in choosing the target signs of original samples that effectively represent their characteristics for image recognition, and then aligning these signals as a proxy task to optimize the condensed dataset. The target signal  $\phi(x;\theta)$  is typically the key information from feature extraction or optimization process using a randomly initialized network parameterized by  $\theta$ . For example, GM aligns the gradient information produced by T with that of the condensed S. Instead, DM matches the embedding distributions of Tand S. As for these approaches in DMF. we can unify the optimization objective as  $\mathcal{L}(S; \theta, T)$ :

$$\mathcal{L}(\mathcal{S};\theta,\mathcal{T}) := \sum_{y\in\mathcal{Y}} \mathcal{D}\big(\mathbb{E}[\phi_{x\sim\mathcal{T}_y}(x;\theta)], \mathbb{E}[\phi_{x\sim\mathcal{S}_y}(x;\theta)]\big), \tag{3}$$

149 where  $\mathbb{E}[\phi_{x \sim \mathcal{T}_y}(x; \theta)] \in \mathbb{R}^C$  and  $\mathbb{E}[\phi_{x \sim \mathcal{S}_y}(x; \theta)] \in \mathbb{R}^C$  are represented expectation vectors of the 150 target signs on  $\mathcal{T}$  and  $\mathcal{S}$ , respectively.  $\mathcal{D}(\cdot, \cdot)$  is a distance function. In DMF, MSE is adopted in DM 151 and DREAM, and MAE is used in IDC. Also, cosine distance is involved in DC.

153 Why does vanilla DD fail to mitigate PA imbalance? Given the PA set  $\mathcal{A}$  in  $\mathcal{T}$ , let us define the 154 class-level sample ratio  $\mathcal{R}_y = \{r_y^{a_1}, r_y^{a_2}, ..., r_y^{a_p}\}$ , where  $r_y^{a_i} = |\mathcal{T}_y^{a_i}|/|\mathcal{T}_y|$ , and  $|\cdot|$  represents the 155 cardinal number of a set. Current DD paradigms focus on preserving TA representativeness for image 156 recognition. Here, we decompose the whole expectation into the expectation of PA-wise groups, i.e., 157  $\mathbb{E}[\phi_{x\sim\mathcal{T}_y}(x;\theta)] = \sum_{a_i\in\mathcal{A}} r_y^{a_i}\mathbb{E}[\phi_{x\sim\mathcal{T}_y^{a_i}}(x;\theta)]$ , and thus Eq. 3 can be rewritten as follows:

$$\mathcal{L}(\mathcal{S};\theta,\mathcal{T}) := \sum_{y\in\mathcal{Y}} \mathcal{D}(\sum_{a_i\in\mathcal{A}} r_y^{a_i} \mathbb{E}[\phi_{x\sim\mathcal{T}_y^{a_i}}(x;\theta)], \mathbb{E}[\phi_{x\sim\mathcal{S}_y}(x;\theta)]).$$
(4)

From Eq. 4, the optimization objective of class y is weighted by the sample ratio  $r_y^{a_i}$  from different groups. When  $\mathcal{T}$  suffers from PA imbalance, e.g.,  $r_y^{a_j} \gg \sum_{i \neq j} r_y^{a_i}$ , the majority group indexed by i contributes more to the alignment compared to minority groups, as present in Fig. 2(a). In other words, S tends to produce more samples belonging to group *i* for the total loss minimization. Therefore, the objective of vanilla DDs suffers from the PA imbalance within T.

Next, we further study how the resulting S is affected by sample ratio  $r_y^{a_i}$  of different groups. To this end, we assume that the optimization process could reach the optimal solution for each class, and as a result, the final resulting S satisfies the condition that the derivative of the objective function with respect to  $\mathbb{E}[\phi_{x\sim S_y}(x;\theta)]$  equals 0. Formally. we have the following mathematical equation:

$$\frac{\partial \mathcal{L}(\mathcal{S}_{y};\theta,\mathcal{T}_{y})}{\partial \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]} = 0 \implies \frac{\partial \mathcal{D}(\sum_{a_{i}\in\mathcal{A}} r_{y}^{a_{i}} \mathbb{E}[\phi_{x\sim\mathcal{T}_{y}}^{a_{i}}(x;\theta)], \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)])}{\partial \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]} = 0 \tag{5}$$

171 Now, let's delve into the specific distance metrics used in vanilla DDs, where the most commonly 172 used metrics are MAE, MSE, and cosine distance. We respectively analyze that the optimal point of 173  $\mathbb{E}[\phi_{x\sim S_y}(x;\theta)]$  could reach under each of these metrics.

176

177 178

169

170 171

 $\frac{\partial \mathcal{L}(\mathcal{S}_{y};\theta,\mathcal{T}_{y})}{\partial \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]} = 0 \implies \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)] = \begin{cases} \sum_{a_{i}\in\mathcal{A}} r_{y}^{a_{i}}\mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], & \text{For MAE} \\ \sum_{a_{i}\in\mathcal{A}} r_{y}^{a_{i}}\mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], & \text{For MSE} \\ \lambda \sum_{a_{i}\in\mathcal{A}} r_{y}^{a_{i}}\mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], & \text{For osine distance} \end{cases}$ (6)

179 180

192

193

212

181 Where  $\lambda$  is a scalar, equaling  $\frac{\|\mathbb{E}[\phi_{x \sim S_y}(x;\theta)]\|_2}{\|\sum_{a_i \in \mathcal{A}} r_y^{a_i} \mathbb{E}[\phi_{x \sim \mathcal{T}_y^{a_i}}(x;\theta)]\|_2}$ . Eq. 6 presents that the expectation of synthetic samples  $\mathbb{E}[\phi_{x \sim S_y}(x;\theta)]$  ultimately converges to an average on expectations of all PA

synthetic samples  $\mathbb{E}[\phi_{x \sim S_y}(x; \theta)]$  ultimately converges to an average on expectations of all PA groups, weighted by their respective sample ratios  $r_y^{a_i}$ . This indicates that vanilla DDs naturally favor majority groups, causing S to shift towards them and inherit their biases.

Therefore, when original datasets suffer from PA imbalance, e.g.,  $r_y^{a_j} \gg \sum_{i \neq j} r_y^{a_i}$ , the unfairness of the synthetic dataset stems from two different aspects: 1) **The majority group renders synthetic samples to locate its region from Eq. 6.** 2) According to Eq. 4, the large sample quantities of the majority group contribute more to the total loss. As a result, **minority groups experience higher loss during testing, which limits the model to represent them accurately.** These factors prompt us to reduce the impact of PA imbalance on the generation of S.

### 4 FairDD

#### 194 4.1 OVERVIEW

195 In this paper, we propose a novel FDD framework that achieves both PA fairness and TA accuracy 196 for the model trained on its generation S, regardless of whether the original datasets exhibit PA 197 fairness. As illustrated in Fig. 2(b), FairDD first partitions the dataset into different groups w.r.t. PA and then introduces an effective synchronized matching to equally align S with each group within  $\mathcal{T}$ . Compared to vanilla DDs, which simply pull the synthetic dataset toward the whole dataset center 199 that is biased toward the majority group in the synthetic dataset, FairDD proposes a group-level 200 synchronized alignment, in which each group attracts the synthetic data toward itself, thus forcing 201 it to move farther from other groups. This "pull-and-push" process prevents the generation from 202 collapsing into majority groups (fairness) and ensures class-level distributional coverage (accuracy). 203

#### 204 4.2 METHODS

As mentioned in Sec. 3, vanilla DDs fail to mitigate PA imbalance and even amplify the discrimination. The relation behind the failure is that the majority group dominates the generation direction of Sand leads to the resulting S inheriting the PA imbalance, i.e., preference to fitting to the majority group. To avoid the synthetic samples collapsing into the majority group, we decompose the single target (dominated by the majority group) into PA-wise sub-targets, and simultaneously align S with these sub-targets, without incorporating the specific sample ratio of each group into the optimization objective. In this way, we obtain the unified objective function of FairDD:

$$\mathcal{L}_{FairDD}(\mathcal{S};\theta,\mathcal{T}) := \sum_{y\in\mathcal{Y}} \sum_{a_i\in\mathcal{A}} \mathcal{D}(\mathbb{E}[\phi_x \,_{\sim\mathcal{T}_y^{a_i}}(x;\theta)], \mathbb{E}[\phi_x \,_{\sim\mathcal{S}_y}(x;\theta)]). \tag{7}$$

The reformulation forms synchronized matching, where different sub-targets attempt to pull S into their corresponding PA regions. Each PA group holds equal importance in generating S, ultimately converging to a balanced (fair) status. Subsequently, we present a theoretical analysis illustrating how FairDD effectively mitigates PA imbalance and guarantees coverage across the entire TA distribution.



Figure 2: Comparison between (a) vanilla DDs and (b) FairDD. Taking C-MNIST (FG) for example 233 (IPC = 10), vanilla DDs directly align S (random initialization) with the whole distribution regardless 234 of majority or minority groups. This causes S to suffer from distributional coverage and inherit 235 (even amplify) the bias of  $\mathcal{T}$ . Instead, FairDD first groups target signals of  $\mathcal{T}$  and then proposes to 236 align  $\mathcal{S}$  (random initialization) with respective group centers. With this synchronized matching,  $\mathcal{S}$  is 237 simultaneously pulled by all group centers in a batch. This prevents the condensed dataset S from 238 being biased towards the majority group, allowing it to better cover the distribution of the original 239 dataset  $\mathcal{T}$ . We can observe that each class in the resulting  $\mathcal{S}$  incorporates multiple colors from the 240 minority color groups and mitigates the bias of  $\mathcal{T}$ , originally dominated by majority colors.

**Theorem 4.1.** For any PA set A, network parameters  $\theta$ , and target signs  $\phi(\cdot)$ ,  $\mathcal{L}_{FairDD}(S; \theta, \mathcal{T})$ could mitigate the influence of PA imbalance of original datasets on generating synthetic samples. Especially when  $\mathcal{D}(\cdot)$  is MAE or MSE, synchronized matching ensures that the signal expectation of S is situated at the center of the expectation across all PA groups within  $\mathcal{T}$ 

246 *Proof.* We assume that  $\mathbb{E}[\phi_{x \sim S_y}(x; \theta)]$  could reach the optimal solution for each class. Consequently, 247 we have  $\partial \mathcal{L}_{FairDD}(S_y; \theta, \mathcal{T}_y) / \partial \mathbb{E}[\phi_{x \sim S_y}(x; \theta)] = 0$ , and derive the respective optimal solution: 248

$$\frac{\partial \mathcal{L}_{FairDD}(\mathcal{S}_{y};\theta,\mathcal{T}_{y})}{\partial \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]} = 0 \implies \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)] = \begin{cases} \frac{1}{|\mathcal{A}|} \sum_{a_{i}\in\mathcal{A}} \mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], & \text{For MAE} \\ \frac{1}{|\mathcal{A}|} \sum_{a_{i}\in\mathcal{A}} \mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], & \text{For MSE} \\ \frac{\lambda}{|\mathcal{A}|} \sum_{a_{i}\in\mathcal{A}} \mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], & \text{For cosine distance} \end{cases}$$
(8)

According to Eq. 8, the resulting  $\mathbb{E}[\phi_{x\sim S_y}(x;\theta)]$  are independent on the sample ratio  $\mathcal{R}_y$ , indicating the corresponding *S* unaffected by  $\mathcal{R}_y$ . As a result, the condensed *S* will not be dominated by majority groups that happened in vanilla DDs. All PA centers contribute equally to the generation of *S*, which succeeds in mitigating the PA imbalance of  $\mathcal{T}$ . Especially when  $\mathcal{D}(\cdot)$  is MAE or MSE, the expectation of target signs of *S* is equal to the arithmetic mean of centers of all PA groups. This shows that *S* generated by FairDD is not biased towards any groups.

Although we mitigate the bias inheritance in vanilla DDs by synchronously aligning S to finegrained PA-wise groups, it is also crucial to investigate whether  $\mathcal{L}_{FairDD}(S;\theta,\mathcal{T})$  (synchronized matching) ensures that the resulting S achieves comprehensive distributional coverage for  $\mathcal{T}$ . As mentioned above,  $\mathcal{L}(S;\theta,\mathcal{T})$  matches S and  $\mathcal{T}$  in a global view to fully cover  $\mathcal{T}$ 's distribution. Below, we provide a theoretical guarantee that  $\mathcal{L}_{FairDD}(S;\theta,\mathcal{T})$  could provide the same or even more comprehensive coverage compared to  $\mathcal{L}(S;\theta,\mathcal{T})$  when  $\mathcal{D}(\cdot,\cdot)$  is a convex distance function, which is commonly used in diverse DDs<sup>1</sup>.

241

249 250

<sup>269</sup> 

<sup>&</sup>lt;sup>1</sup>Emperical experiments show FairDD also can cover the TA distributions when  $\mathcal{D}(\cdot, \cdot)$  is not convex.

**Theorem 4.2.** For any PA set  $\mathcal{A}$  and target signs  $\phi_{\theta}(\cdot)$ ,  $\mathcal{L}_{FairDD}(S; \theta, \mathcal{T})$  is the upper bound of vanilla unified objective  $\mathcal{L}(S; \theta, \mathcal{T})$ , i.e.,  $\mathcal{L}_{FairDD}(S; \theta, \mathcal{T}) \geq \mathcal{L}(S; \theta, \mathcal{T})$ , when  $\mathcal{D}(\cdot, \cdot)$  is convex. Optimizing  $\mathcal{L}_{FairDD}(S; \theta, \mathcal{T})$  can guarantee the comprehensive distribution coverage for  $\mathcal{T}$ . Proof.

$$\mathcal{L}(\mathcal{S};\theta,\mathcal{T}) = \sum_{y\in\mathcal{Y}} \mathcal{D}\left(\mathbb{E}[\phi_{x\sim\mathcal{T}_{y}}(x;\theta)], \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]\right)$$
$$= \sum_{y\in\mathcal{Y}} \mathcal{D}\left(\sum_{a_{i}\in\mathcal{A}} r_{y}^{a_{i}} \mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]\right)$$
$$\leq \sum_{y\in\mathcal{Y}} \sum_{a_{i}\in\mathcal{A}} r_{y}^{a_{i}} \mathcal{D}\left(\mathbb{E}[\phi_{x\sim\mathcal{T}_{y}^{a_{i}}}(x;\theta)], \mathbb{E}[\phi_{x\sim\mathcal{S}_{y}}(x;\theta)]\right)$$
(9)

$$\leq \sum_{y \in \mathcal{Y}} \sum_{a_i \in \mathcal{A}} \mathcal{D} \left( \mathbb{E}[\phi_{x \sim \mathcal{T}_y^{a_i}}(x; \theta)], \mathbb{E}[\phi_{x \sim \mathcal{S}_y}(x; \theta)] \right)$$
(10)

$$=\mathcal{L}_{FairDD}(\mathcal{S};\theta,\mathcal{T})$$

Eq. 9 is obtained according to the Jensen Inequality, and Eq. 10 is given since group ratios are smaller than one.  $\mathcal{L}_{FairDD}(S; \theta, T)$  serves as the upper bound of  $\mathcal{L}(S; \theta, T)$ , meaning that minimizing  $\mathcal{L}_{FairDD}(S; \theta, T)$  ensures the minimization of  $\mathcal{L}(S; \theta, T)$ . Hence, optimizing S with FairDD can guarantee the distributional coverage, at least as comprehensive as  $\mathcal{L}(S; \theta, T)$  tailored for accuracy.

#### 289 5 EXPERIMENT

281

283 284

285

286

287

288

#### 290 5.1 EXPERIMENT SETUP

**Datasets** Comprehensive experiments have been conducted on publicly available datasets of diverse 292 biases, including foreground bias (FG), background bias (BG), BG & FG bias, and real-world bias. 293 C-MNIST (FG) is a variant of MNIST (LeCun et al., 2010) used to evaluate model fairness, where the 294 handwriting numbers in each class are painted with ten different colors. To correlate the TA (digital 295 number) and PA (color) within the training dataset, each training class is predominantly associated 296 with one color according to the same biased ratio (BR), while the remaining samples are evenly 297 painted with the other nine colors. BR is the ratio of the majority group samples to the total samples 298 across all groups. For the test dataset, we evenly paint the numbers for each class with ten colors to 299 test the model bias trained on S. C-MNIST (BG) adopts the same operation on the background and 300 keeps the foreground unchanged. Colored-FMNIST (FG) is the modified version of Fashion-MNIST, originally aiming to classify object semantics. Like C-MNIST (FG), we color the objects for the 301 training and test datasets. Colored-FMNIST (BG) paints the background similarly to C-MNIST 302 (BG). CIFAR10-S (BG & FG) introduces a PA by applying grayscale or not to CIFAR10 samples. 303 Following Wang et al. (2020), we grayscale a portion of the training images, correlating TA and PA 304 among different classes. For fairness evaluation, we duplicate the test images, apply grayscale to the 305 copies, and add them to the test dataset. We also test FairDD on the real-world facial dataset CelebA, 306 a widely used fairness dataset. We follow the common practice of treating attractive attribute as 307 TA and gender as PA (evaluations on other attributes refer to Appendix B). 308

**Baselines & Evaluation metrics** FairDD is a general fairness framework applicable to diverse DDs in DMF. We apply FairDD to diverse DMF approaches including DM method DM (Zhao & Bilen, 2023) and GM methods DC (Zhao et al., 2020), IDC (DC version) (Kim et al., 2022), and DREAM (DC version) (Liu et al., 2023). To provide an overall evaluation for model bias toward PA, we use DEO<sub>M</sub>( $\downarrow$ )  $\in$  [0, 100] and DEO<sub>A</sub>( $\downarrow$ )  $\in$  [0, 100] to measure the worst and average fairness levels. Also, we report accuracy( $\uparrow$ ) to assess the model's prediction of TA. We also provide a comparison with MTT in Appendix L. Sometimes, we will abuse DM+FairDD and FairDD for clarification.

Implementation details We default to BR of 0.9 for all synthetic original datasets to induce significant PA skew. In Table 24, we conduct the ablation study on BR. All baselines are reproduced using official implementations. FairDD doesn't introduce extra hyperparameters or learnable parameters. Experiments are conducted on PyTorch 2.0.0 with a single NVIDIA RTX 3090 24GB GPU.

320 5.2 MAIN RESULTS

We use distilled datasets S from different DDs to train and evaluate ConvNet with the same parameters, and then report the corresponding fairness and accuracy. *Random* refers to sampling defined IPC from the original dataset to create smaller datasets. Besides, *Whole* means we train the model using the entire training dataset without distillation or sampling.

	Table 1	: Fair	ness	com	paris	son o	on di	vers	e IPO	Cs. T	The b	est r	esult	s are	hig	hligh	nted i	in <b>b</b> o	ld.	
Methods Dataset	IPC R	andom M DEO <sub>A</sub>	DEO <sub>M</sub>	M DEO <sub>A</sub> I	DM+F DEO <sub>M</sub>	airDD DEOA	D DEO <sub>M</sub>	C DEO <sub>A</sub>	DC+F DEO <sub>M</sub>	TairDD	IE DEO <sub>M</sub>	DEO <sub>A</sub>	IDC+F DEOM	airDD DEO A	DRE DEO <sub>M</sub>	$\frac{AM}{DEO_A}$	+Fai DEOM	$\frac{rDD}{DEO_A}$	Wh DEO <sub>M</sub>	ole DEO <sub>A</sub>
C-MNIS (FG)	Г 10 100	.0 98.72 .0 99.58	100.0	99.96 91.68	17.04 10.05	7.95 5.46	99.85	65.61 20.55	26.75	11.96 8.86	100.0	91.45 34.91	12.24 9.18	6.64 5.94	98.99 52.03	78.71 26.63	11.88 18.37	7.21	10.10	5.89
C-MNIS	T   10   100 T   50   100	.0 88.64 .0 99.11 .0 99.77	99.36	99.97 97.85	8.17 13.42 8.98	4.86 6.77 5.25	45.27 100.0 60.66	73.60 26.38	22.32 20.66 20.29	9.49 9.94 9.90	100.0 93.05	35.82 88.30 42.23	11.88 18.61 19.66	6.21 7.50 8.05	69.30 100.0 64.15	52.06 23.30	11.88 15.31 20.41	6.88 6.83 9.04	9.70	5.78
(BC)	100 100	.0 89.07	100.0	52.23	6.60	4.31	62.63	20.87	32.58	10.40	63.24	27.79	12.24	6.32	44.88	22.86	16.33	7.80		
C-FMNIS (FG)	T 10 100 50 100 100 100	.0 99.18 .0 94.61 .0 94.85	100.0 100.0 100.0	99.05 96.46 85.11	26.87 24.92 23.83	16.38 13.74 12.75	99.40 99.33 99.58	78.96 67.02 66.45	46.80 46.67 56.68	24.01 21.48 23.07	100.0 100.0 100.0	97.27 81.93 79.10	32.33 40.00 48.33	16.80 17.37 17.43	100.0 99.67 97.33	95.17 83.27 70.10	42.00 47.67 74.00	20.87 22.33 40.40	79.20	41.72
C-FMNIS (BG)	T 10 100 50 100 100 100	.0 99.40 .0 98.52 .0 96.05	100.0 100.0 100.0	99.68 99.71 93.88	33.05 24.50 21.95	19.72 14.47 13.33	100.0 100.0 99.70	92.91 75.41 73.38	61.75 44.60 52.75	34.88 25.25 23.48	100.0 100.0 100.0	99.40 95.60 90.70	42.00 78.00 77.00	23.80 34.50 36.00	100.0 100.0 100.0	94.70 88.40 83.90	36.00 34.00 40.00	23.50 23.70 23.20	91.40	51.68
CIFAR10	-S 50 57.1 100 66.4	04 8.29 11 28.89 19 43.16	59.20 75.13 73.81	39.31 55.70 55.10	31.75 18.28 14.77	8.73 7.35 5.89	42.23 71.46 68.69	27.35 45.81 48.64	22.08 34.39 32.70	8.22 11.21 11.26	80.70 92.00 92.70	48.38 60.56 60.93	19.90 29.00 62.80	5.28 9.10 25.18	51.80 56.80 82.30	31.43 36.19 48.12	20.80 14.70 12.10	7.77 6.53 6.06	49.72	33.17
CelebA	10 10.4 50 22.8 100 18.6	48 9.20 38 20.32 57 18.01	30.01 40.26 42.63	28.85 38.81 41.12	9.37 14.08 10.93	5.71 9.87 6.65	15.48 24.89 29.00	14.16 23.83 27.52	6.64 14.33 18.16	5.29 12.92 17.04	34.85 56.74 50.99	34.48 46.50 42.66	8.36 22.57 28.27	4.49 15.15 17.63	40.75 43.57 52.51	36.70 38.53 39.34	9.20 23.62 24.87	5.36 14.29 15.36	24.85	24.16
000	0000	000	00	00	000		0											1	1	
222	2 2 2 2 2	2222	22	22	2 2 2	222	2			Å Å 📒						A				2 44
444	9 <b>4</b> 4 4	444	44	5 <b>5</b> 4 9	4 4 4	44	4													
555	5555	5 5 5 5 5 6 6 6	55	55	5 C	5 5 6 6 6	6	* <i>\$</i> 4 .} 												
777		771	77	77	77	777	7		-seles 🎿			nititi antia <b>anti</b> Stati	ي هي <u>هر</u> و 200 ون	2 <mark>2 - 22</mark> 22 - 22 - 22 23 - 22						100 A
999	999	999	99	99	99	999	9	ار نار پ	11		4.4	1 13	A.	A A J		-16			6	88 <b>9</b>
	DM				DM				DI	И				JМ		(A)		DN		
000	0000	1/1	00	00	000 11	000	0				1 N						<b>8</b> 8			
2 $2$ $23$ $3$ $2$	2 2 2 3	2222	22	22	a 2 .	222	2	N 🕺 ří		AA		A 8								2 2
444	4 4 4 4	1444	4 4	44	44	4 4	4	1	Ä 👗	🦲 🛄 🧧	à 🛕	88		8 <b>8</b> 8	<u>8</u>	<u>.</u>				-
555	558	5555 6666	5 5 6 6	55 66	5 5	555	6	* <i></i> ?   🔒 👭	1992 - 1992 - 	aa 🛞 🔒 🏅				2 2 () ()	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1					
777	7777	777	77	77	777	777	7	ia è		مر میں ایک ایک ایک 🛄	- 🦽 🛷		en e	a <mark>as</mark> aa Ngara	0 - 200 <mark>- 200</mark> - 220 - 200					
999	9999	999	99	99	99	999	9	- <u></u>		<u>, ,</u>		A di	<u> </u>	* * 3	<u>.</u>					
Fa	irDD+E	DМ	ł	FairD	D+I	DM		Fa	irDD	+DN	1	F	FairD	D+D	M		Fair	DD-	-DM	
(a) C-	MNIST	(FG).	(b)	C-M	NIST	Г (ВС	3). (0	c) C-	FMN	IST	(FG).	(d) (	C-FM	INIS	Г (ВС	3).	(e) C	CIFA	R10-3	S.
Figure	e 3: Vi	sualiz	ation	con	npar	ison	on	${\cal S}$ at	IPC	2 = 1	0 for	r div	erse	data	sets.	Fai	irDD	) suc	cess	fully
mitiga	ites the	bias f	rom	orig	inal	data	sets	in (a	l) (fo	regr	ound	digi	tal c	olor	), (b)	bac	kgro	und	colo	r, (c)
foregr	ound o	bject o	color,	, (d)	back	grou	and o	coloi	r, and	d (e)	fore	grou	nd ai	nd ba	ackg	roun	d gra	iysca	ıle.	
FairD	D sign	ificant	tly in	ipro	ves t	he fa	airn	ess o	f vaı	nilla	DD a	appr	oach	ies	We	prov	ide c	omp	rehe	nsive
fairnes	ss comp	parisoi		ross '	varic	ous L	DD p	arad	igms	, inc	ludin	ig Df	VI an	d D(	J. As	s illus	strate	ed in	Tabl	le 21,
towar	t DDS I	ed gr	nnuş	Jale	$C_{-N}$	JIAS J	ST (	FG)	n uic the	diet	filled	uata dat	SELS acete	anu fro	even m D	M s	uffer	· fro	m se	Nere
unfair	ness at	· IPC=	-10 c	omr	oared		Who	le v	, inc vith	DEC	) A a	nd F	DEO.	res	in D	σ 1(	$\frac{1101}{0}$	and	00 Q	6 vs
10.10	and $5.8$	39. In	some	case	es. R	and	om t	rese	nts b	better	fair	ness	than	van	illa I	DDs.	part	icula	rlv v	when
dealin	g with	compl	ex ob	oject	s like	e Cel	lebÅ	. Th	is su	gges	ts tha	it wh	ile v	anill	a DI	)s ef	fecti	vely	cond	lense
inforn	nation	into si	malle	er sa	mple	es, th	neir i	indu	ctive	bias	s, wł	nich	favo	rs th	e ma	ijori	ty gr	oup,	woi	sens
the fa	irness 1	to the	mino	ority	gro	up.	How	veve	r, wł	nen I	FairE	DD is	s app	olied	to v	anil	la D	Ds, 1	there	e is a
signifi	cant in	nprove	emen	t in f	airn	ess p	berfo	orma	nce,	with	DE	$D_M d$	lropp	oing	subs	tanti	ally	from	100	0.0 to
17.04,	and D	EO <sub>A</sub> (	decre	asin	g fro	pm 9	19.96	to '	1.95	in C	-MN	IST area	(FG	i). Ti	115 11	1d1C8	tes 1	that	Fairl	JD'S
synchi		match	ung ( rhate	EIISU:	res II		jual ber	ueat	men	t Of E	acn j	grou rigin	p, en	prec	vely	1111[18 n th	gaun	g ine	: Dias	s inat
For ev	a DDS ( ample	$DC \perp$	Fairl	. rai	עעי חוונות	rfor	me V	vhol	e in l	LIC D	148 0 ANIS	ngin T (F	iany 'G) a	nd C	TFA1	ui¢ R10-	S as	gilla Wel	i uata Laci	asets. n the
real-w	orld da	taset (	Celeb	A. a	chie	ving	the 4	over	all in	nprov	veme	nt or	טע ו DF	III C	and 1	DEO	$D_{\Lambda}$ m	etric	5. Si	milar
perfor	mance	gains	are a	also (	obse	rved	linc	other	base	eline	s. Tł	nese	resu	lts u	nder	score	e the	effe	ctive	eness
and ve	ersatilit	y of F	airDl	D. W	le pr	ovid	e the	e vis	ualiz	atior	n con	npari	ison	of re	sulti	ng S	in F	Fig. 3	3.	

FairDD maintains the comparable and even higher accuracy than vanilla DD approaches
 A fairness framework must maintain TA accuracy in addition to improving fairness across PA groups.
 We report the TA accuracy of FairDD in comparison to other baselines in Table 22. Compared to

378 379

Table 2: Accuracy comparison	on div	verse IPCs.
------------------------------	--------	-------------

T 11	0	<b>C</b> 1	•
Table	- <b>X</b> •	( ross-arch	comparison
Iuoio	<i>J</i> •	Crobb uren.	companyon.

						1						
Methods Datasets	IPC	Random Acc.	$\left  \frac{\mathrm{DM}}{\mathrm{Acc.}} \right $	+FairDD Acc.	$\left  \frac{DC}{Acc.} \right $	+FairDD Acc.	IDC Acc.	+FairDD Acc.	DREAM	+FairDD Acc.	Whole Acc.	Method
C-MNIST (FG)	10 50 100	30.75 47.38 67.41	25.01 56.84 78.04	94.61 96.58 96.79	71.41 90.54 91.64	90.62 92.68 93.23	53.06 88.55 90.39	95.67 96.77 97.11	75.04 91.02 88.87	94.04 94.59 95.16	97.71	C-MNIST (FG)
C-MNIST (BG)	10   50   100	27.95 45.52 67.28	23.40 47.74 79.87	94.88 96.86 97.33	65.91 88.53 90.20	90.84 92.20 92.73	62.09 86.14 89.66	94.84 95.29 95.84	79.81 89.24 90.70	93.54 93.20 94.06	97.80	C FADUET
C-FMNIST (FG)	10   50   100	32.80 42.48 55.31	33.35 49.94 57.99	77.09 82.11 83.25	60.77 69.08 68.84	76.01 75.83 74.91	44.08 64.45 66.37	79.66 80.80 80.28	49.72 65.69 68.25	77.24 78.79 78.51	82.94	(BG)
C-FMNIST (BG)	10   50   100	24.96 34.92 44.87	22.26 36.27 49.30	71.10 79.07 80.63	47.32 60.58 62.70	68.51 75.80 71.76	37.59 46.20 48.61	72.67 73.72 73.18	45.30 53.62 53.32	71.56 72.80 73.00	77.97	CIFAR10-S
CIFAR10-S	10   50   100	23.60 36.46 39.34	37.88 45.02 48.11	45.17 58.84 61.33	37.88 41.28 42.73	41.82 49.26 51.74	48.30 47.26 47.27	56.40 57.84 56.98	55.09 57.59 57.14	58.40 61.85 62.70	69.78	
CelebA	10   50   100	54.51 55.99 60.62	61.79 64.61 65.13	64.37 68.50 68.84	57.19 60.16 62.53	<b>57.63</b> 59.89 61.89	61.49 60.75 64.04	63.54 66.89 67.24	64.38 64.62 62.58	66.26 68.26 64.12	74.09	CelebA

Mathod	Cross		DM		DN	1+FairD	D
Wiethou	arch.	DEOM	DEOA	Acc.	DEOM	DEOA	Acc
	ConvNet	100.0	91.68	56.84	10.05	5.46	96.5
C MNIST	AlexNet	100.0	98.82	44.02	10.35	6.16	96.1
(EC)	VGG11	99.70	70.73	75.22	9.55	5.39	96.8
(FG)	ResNet18	100.0	96.00	52.05	8.40	4.63	97.1
	Mean	99.93	89.31	57.03	9.59	5.41	96.6
C-FMNIST (BG)	ConvNet	100.0	99.71	36.27	24.50	14.47	79.0
	AlexNet	100.0	99.75	22.72	20.60	14.11	76.1
	VGG11	100.0	97.77	43.11	21.60	14.36	78.5
	ResNet18	100.0	99.78	23.37	22.50	14.96	75.2
	Mean	100.0	99.25	31.37	22.30	14.73	77.2
	ConvNet	75.13	55.70	45.02	18.28	7.35	58.8
	AlexNet	75.30	52.57	36.09	15.84	5.12	49.1
CIFAR10-S	VGG11	61.48	44.05	43.23	11.51	4.16	52.0
	ResNet18	76.23	54.35	38.03	16.44	5.14	50.9
	Mean	72.04	51.67	40.59	15.27	5.44	52.9
	ConvNet	40.26	38.81	64.61	14.08	9.87	68.
	AlexNet	32.51	31.62	63.10	9.38	5.75	64.2
CelebA	VGG11	26.03	24.63	61.57	8.95	6.32	62.0
	ResNet18	25.60	24.93	60.32	6.72	4.29	61.8
	Mean	31.10	30.25	62.40	9.78	6.58	64.1

394 *Random*, training the model by vanilla DDs yields better performance. This shows that vanilla DDs capture the informative patterns of majority groups, improving their TA accuracy. However, by 396 focusing on dominant patterns in majority groups, they neglect the important patterns in minority 397 groups within the training datasets. Thus, their representation coverage is limited. In contrast, FairDD proposes synchronized matching to push the S to cover each group, and as a result, the generated S398 retains key patterns of all groups and achieves comprehensive coverage. For example, DM obtains 399 25.01 at IPC = 10 on C-MNIST (FG), and its accuracy boosts to 94.61 when applying FairDD. In 400 real-world CelebA, FairDD obtains comparable performance for DC and presents superiority over 401 vanilla DDs. These demonstrate that FairDD could mitigate the bias without compromising accuracy. 402

**Generalization to diverse architectures** Here, we investigate the cross-model generalization of 403 FairDD, where ConvNet is used to condense datasets, and we evaluate S on other architectures, 404 including AlexNet, VGG11, and ResNet18. We compare DM and FairDD across four datasets at 405 IPC = 50, evaluating performance against BG, FG, BG & FG, and real-world biases. As shown in 406 Table 23, among these architectures, FairDD achieves DEO<sub>M</sub> of 10.05, 10.35, 9.55, and 8.40 on 407 C-MNIST (FG), DEO<sub>A</sub> of 14.47, 14.11, 14.36, and 14.96 on C-FMNIST (BG), and accuracy of 58.84, 408 49.16, 52.65, and 50.93 on CIFAR10-S. These steady results suggest that S generated by FairDD 409 are not restricted to the model used for distillation but generalize well across diverse architectures. 410 Additionally, with the model capacity increasing, the model generally tends to be more fair to all 411 groups. However, the accuracy sometimes decreases, such as when it drops from 58.84 (ConvNet) 412 to 50.93 (ResNet18) in CIFAR10-S and from 68.50 (ConvNet) to 61.80 (ResNet18) in CelebA. We assume that while increased attention from larger models can lead to accuracy gains for minority 413 groups, it may limit the representations for majority groups at certain levels. The accuracy gains for 414 minority groups may be smaller than the accuracy losses for majority groups, particularly in larger 415 models that have limited potential improvement in recognizing minority groups. As a result, overall 416 accuracy may decrease despite fairness improvement. 417

418 5.3 RESULT ANALYSIS

Visualization analysis on repre-419 sentation coverage We investigate 420 whether the FairDD effectively cov-421 ers the whole distribution of original 422 datasets. For this purpose, we first 423 feed the original training set into the 424 randomly initialized network used in 425 the distillation, to extract the corre-426 sponding features. Subsequently, we 427 use the same network to extract fea-428 tures of distilled dataset S from DM and FairDD. As shown in Fig. 4(a), 429 the synthetic samples in vanilla DDs 430 almost locate the majority group for 431



(a) The  ${\mathcal S}$  distribution of generated by DM.

(b) The  $\mathcal S$  distribution of generated by FairDD.

Figure 4: Feature coverage comparison on TA between DM and DM+FairDD. The training and synthetic dataset features are extracted by  $\phi_{\theta}$ . We visualize one class and make the remaining classes transparent for presentation. The *S* generated by DM and FairDD are marked by stars in (a) and (b).

optimizing the original alignment objective. In this case, vanilla DDs neglect to condense the key



Figure 5: T-SNE visualization comparison on PA and TA between DM and DM+FairDD. Each color represents distinct PA groups in (a) and (b). (c) and (d) use different colors to represent different PA. In (a), DM shows obvious distinctiveness towards different PA, whereas in (b), DM+FairDD eliminates the recognition of PA. As illustrated in (c) and (d), DM+FairDD optimizes the model to produce compact TA representations, but DM fails to cluster TA from the same class.

Table 4	: Ablati	on or	n BR	t at I	PC =	= 50.	able 5: A	blatio	n on	initia	ıliza	tion	at IP	C =
Methods Dataset	BR $\left  \frac{1}{\text{DEO}_{M}} \right $	DM DEO <sub>A</sub>	Acc.	DN DEO <sub>M</sub>	A+FairI DEO <sub>A</sub>	DD Acc.	Methods Dataset	Init.	DEOM	$\frac{\text{DM}}{\text{DEO}_{\text{A}}}$	Acc.	DN DEO <sub>M</sub>	4+FairD DEO <sub>A</sub>	D Acc.
C-MNIST (FG)	0.85 99.54 0.90 100.0 0.95 100.0	70.13 91.68 100.0	76.24 56.84 33.73	10.13 10.05 10.30	5.20 5.46 5.84	96.62 96.58 96.05	C-MNIST (FG)	Random Noise Hybrid	100.0 100.0 100.0	91.68 99.64 99.06	56.84 41.45 39.97	10.05 9.33 9.03	5.46 5.28 5.33	96.58 96.06 96.27
C-FMNIST (BG)	0.85100.00.90100.00.95100.0	95.54 99.71 99.79	46.14 36.27 26.30	23.75 24.50 29.15	13.85 14.47 17.72	79.61 79.07 78.46	C-FMNIST (BG)	Random Noise Hybrid	100.0 100.0 100.0	99.71 99.67 99.68	36.27 22.92 26.38	24.50 23.00 21.45	14.47 14.40 14.34	79.07 78.84 79.19
CIFAR10-S	0.85 71.75 0.90 75.13 0.95 75.43	50.11 55.70 58.58	46.99 45.02 43.56	16.44 18.28 17.49	6.58 7.35 7.10	59.12 58.84 58.18	CIFAR10-S	Random Noise Hybrid	75.13 55.28 65.59	55.70 37.26 46.18	45.02 46.97 45.16	18.28 16.15 17.30	7.35 6.13 6.71	58.84 56.41 56.78

patterns of minority groups. This leads to the information loss of minority groups in S. FairDD achieves overall coverage for both majority and minority groups in Fig. 4(b). This is because FairDD introduces synchronized matching to reformulate the distillation objective for aligning the PA-wise groups rather than being dominated by the majority group like vanilla DDs. In doing so, FairDD avoids S collapsing into the majority group and retains informative patterns from all groups.

Visualization analysis on fairness and accuracy To intuitively present the effectiveness of 460 FairDD, we train  $g_{\psi}$  using S of C-MNIST (FG) distilled by DM and FairDD, and then extract the 461 features from the test dataset. Different colors paint these resulting features according to PA and 462 TA, respectively. As shown in Figs. 5(a) and 5(b), features with the same PA tend to form a cluster, 463 indicating that the model trained on DM is sensitive to PA and thus failing to guarantee fairness 464 among all PA. In contrast with DM, the feature distributions in Fig. 5(b) exhibit nearly complete 465 overlaps across all PA. It shows that the model trained on FairDD is agnostic to PA and does not 466 exhibit bias towards these PA. Besides the PA fairness, we also study the feature distribution from the 467 TA perspective. Fig. 5(c) shows that features belonging to one TA scatter and fail to provide compact 468 representations for one class. The failure of DM can be attributed to model bias toward PA. Combined 469 with Fig. 5(a), it can be observed that PA has a stronger influence on the feature distribution compared 470 to TA. As a result, PA-wise representations are tightly clustered, but representations from the same TA are divided into PA-wise parts. In contrast, FairDD proposes synchronized matching effectively 471 mitigates this by treating each PA group equally within one TA. The equal treatment allows different 472 PA groups within the same TA to cluster more easily, leading to more compact representations that 473 benefit capturing class semantics in Fig. 5(d). These results highlight the superiority of FairDD in 474 improving PA fairness and TA accuracy. We also provide visualization analysis on S generation in 475 Appendix M. Additional computation overhead is provided in Appendix C. 476

477 5.4 ABLATION STUDY

440

441

442

443

444 445

478 Ablation on biased ratio of original datasets BR reflects the extent of unfairness in the original 479 datasets and indicates the level of PA skew that the distillation process of S will encounter. We 480 investigate the impact of BR values on fairness performance by setting BR to  $\{0.85, 0.90, 0.95\}$  on 481 C-MNIST (FG), C-FMNIST (BG), and CIFAR10-S. The results at IPC = 50 in Table 24 show that 482 DM is sensitive to the BR of original datasets, with its  $DEO_M$  decreasing from 70.13 to 100.0 as BR increases from 0.85 to 0.95. A similar trend is observed in other datasets. Compared to DM, 483 FairDD maintains consistent fairness and accuracy levels across different biases. This is attributed to 484 the synchronized matching, which explicitly aligns each PA-wise subtarget, reducing sensitivity to 485 group-specific sample numbers. This shows FairDD's robustness to PA skew in the original datasets.

486 Ablation on initialization of synthetic images The initialization of S determines the prior in-487 formation obtained by DDs. We examine the effect of different initialization using three strategies: 488 random: randomly drawing samples from the original datasets to initialize S; *Noise*: using noise 489 obeying the standard normal distribution for initialization; and hybrid: selecting each initialization 490 sample with equal probability to follow either the random or Noise strategy. In Table 25, the  $DEO_M$ and DEOA metrics of DM have largely fluctuate under these initialization strategies. This suggests 491 that DM fails to incorporate sufficient distilled information into  $\mathcal{S}$ , making it overly dependent on 492 the prior information of S. In contrast, FairDD achieves robust performance in both fairness and 493 accuracy, indicating that the informative patterns during distillation dominate the generation of  $\mathcal{S}$ . 494

## <sup>495</sup> 6 Related Work

496 **Dataset distillation** Dataset distillation has been broadly applied to many important fields (Lee 497 et al., 2023; He et al., 2024; Feng et al., 2023; Chen et al., 2023). The first work (Wang et al., 498 2018) attempts to formulate dataset distillation as a bi-level optimization problem. However, the two 499 folds of the optimization process are time-consuming. Neural tangent kernel (Jacot et al., 2018) are 500 utilized to obtain the close form of the inner loop (Nguyen et al., 2021; Loo et al., 2022; Zhou et al., 501 2022). Some works propose surrogate objectives to achieve comparable even better performance, 502 including matching-based methods like GM (Zhao et al., 2020; Zhao & Bilen, 2021; Lee et al., 503 2022a), DM (Zhao & Bilen, 2023; Wang et al., 2022; Zhao & Bilen, 2022), TM (Cazenavette 504 et al., 2022; Cui et al., 2022), soft label learning (Bohdal et al., 2020; Sucholutsky & Schonlau, 2021) and factorization (Kim et al., 2022; Deng & Russakovsky, 2022; Liu et al., 2022; Lee et al., 505 2022a). Cui et al. (2024) and Zhao et al. (2024) focus on reducing the bias to improve the classification 506 performance (Vahidian et al., 2024; Wang et al., 2024). In summary, current DD approaches only 507 pursue classification accuracy that models trained on synthetic datasets, while they neglect fairness 508 concerning PA. Therefore, we propose FairDD to improve fairness without sacrificing TA accuracy. 509

**Visual fairness** With the advancement of computer vision, fair predictions without discrimination 510 towards minority groups have become a crucial topic (Caton & Haas, 2024). According to the stage 511 of bias mitigation, the research field of fairness algorithm (Bellamy et al., 2019) can be classified into 512 three branches: Pre-processing (Creager et al., 2019; Louizos et al., 2015; Quadrianto et al., 2019; 513 Sattigeri et al., 2019), In-processing (Agarwal et al., 2018; Jiang & Nachum, 2020; Zafar et al., 2017; 514 Zhang et al., 2018; Jung et al., 2021; Wang et al., 2020; Jung et al., 2022; Zemel et al., 2013), and 515 Post-processing (Alghamdi et al., 2020; Hardt et al., 2016). Our research falls within Pre-processing 516 branch. Pre-processing aims to generate a fair version of datasets for downstream tasks. Several 517 related works frame this issue as a data-to-data translation problem, utilizing generative models to 518 produce fairer datasets concerning protected groups (Sattigeri et al., 2019; Quadrianto et al., 2019). 519 However, unlike the traditional fairness approach that primarily focuses on reducing bias in original 520 datasets Subramanian et al. (2021); Tarzanagh et al. (2023); Vogel et al. (2020); Rangwani et al. (2022), our work aims to integrate fairness into DD. Our objective is to mitigate the bias of original 521 datasets while simultaneously condensing them into smaller and more informative counterparts. 522

<sup>523</sup> 7 CONCLUSION

524 This paper reveals for the first time that vanilla DDs fail to mitigate the bias of original datasets and 525 even exacerbate the bias. To address the problem, we propose a unified fair dataset distillation frame-526 work called FairDD, broadly applicable to various DDs in DMF. FairDD requires no modifications 527 to the architectures of vanilla DDs and introduces an easy-to-implement yet effective attribute-wise 528 matching. This method mitigates the dominance of the majority group and ensures that synthetic 529 datasets equally incorporate representative patterns with all protected attributes from both majority 530 groups and minority groups. By doing so, FairDD guarantees the fairness of synthetic datasets while maintaining their representativeness for image recognition. We provide extensive theoretical analysis 531 and empirical results to demonstrate the superiority of FairDD. 532

Limitation Since FairDD relies on PA's prior information to conduct attribute-wise matching, it is
 valuable to explore the scenario where PA is unavailable Liu et al. (2021). A potential solution is to
 generate pseudo-labels to guide FairDD through self-supervised learning or unsupervised learning.

Board impact This paper aims to improve data efficiency and enhance model fairness in modern
machine learning, fully compliant with legal regulations. Since training a fair model from scripts
with extensive data is time-consuming, our work in providing a fair condensed dataset for effective
model training can have significant societal impacts. We hope our research raises attention to achieve
both fairness and accuracy for dataset distillation in academia and industry.

## 540 REFERENCES 541

542 543 544	Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. In <i>International Conference on Machine Learning</i> , pp. 60–69. PMLR, 2018.
545 546 547	Wael Alghamdi, Shahab Asoodeh, Hao Wang, Flavio P Calmon, Dennis Wei, and Karthikeyan Nate- san Ramamurthy. Model projection: Theory and applications to fair machine learning. In 2020 IEEE International Symposium on Information Theory (ISIT), pp. 2711–2716. IEEE, 2020.
548 549 550 551 552	Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilović, et al. Ai fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. <i>IBM Journal of Research and Development</i> , 63(4/5):4–1, 2019.
553 554	Ondrej Bohdal, Yongxin Yang, and Timothy Hospedales. Flexible dataset distillation: Learn labels instead of images. <i>arXiv preprint arXiv:2006.08572</i> , 2020.
555 556 557	Simon Caton and Christian Haas. Fairness in machine learning: A survey. <i>ACM Computing Surveys</i> , 56(7):1–38, 2024.
558 559 560	George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Dataset distillation by matching training trajectories. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 4750–4759, 2022.
561 562 563	Xuxi Chen, Yu Yang, Zhangyang Wang, and Baharan Mirzasoleiman. Data distillation can be like vodka: Distilling more times for better quality. <i>arXiv preprint arXiv:2310.06982</i> , 2023.
564 565 566	Ming-Yu Chung, Sheng-Yen Chou, Chia-Mu Yu, Pin-Yu Chen, Sy-Yen Kuo, and Tsung-Yi Ho. Rethinking backdoor attacks on dataset distillation: A kernel method perspective. <i>arXiv preprint arXiv:2311.16646</i> , 2023.
567 568 569	Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In <i>International conference on machine learning</i> , pp. 1436–1445. PMLR, 2019.
570 571 572	Justin Cui, Ruochen Wang, Si Si, and Cho-Jui Hsieh. Scaling up dataset distillation to imagenet-1k with constant memory. <i>arXiv preprint arXiv:2211.10586</i> , 2022.
573 574	Justin Cui, Ruochen Wang, Yuanhao Xiong, and Cho-Jui Hsieh. Ameliorate spurious correlations in dataset condensation. <i>arXiv preprint arXiv:2406.06609</i> , 2024.
575 576 577	Zhiwei Deng and Olga Russakovsky. Remember the past: Distilling datasets into addressable memories for neural networks. <i>arXiv preprint arXiv:2206.02916</i> , 2022.
578 579 580 581 582	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. <i>ArXiv</i> , abs/2010.11929, 2020. URL https://api.semanticscholar.org/CorpusID: 225039882.
583 584 585	Yunzhen Feng, Shanmukha Ramakrishna Vedantam, and Julia Kempe. Embarrassingly simple dataset distillation. In <i>The Twelfth International Conference on Learning Representations</i> , 2023.
586 587 588	Ziyao Guo, Kai Wang, George Cazenavette, Hui Li, Kaipeng Zhang, and Yang You. Towards lossless dataset distillation via difficulty-aligned trajectory matching. <i>arXiv preprint arXiv:2310.05773</i> , 2023.
589 590 591	Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. <i>Advances in neural information processing systems</i> , 29, 2016.
592 593	Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778, 2015. URL https://api.semanticscholar.org/CorpusID:206594692.

594 595 596	Yang He and Joey Tianyi Zhou. Data-independent module-aware pruning for hierarchical vision transformers. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=7016foUi1G.
597 598 599 600	Yang He, Lingao Xiao, Joey Tianyi Zhou, and Ivor Tsang. Multisize dataset condensation. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=FVhmnvqnsI.
601 602 603	Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. <i>Advances in neural information processing systems</i> , 31, 2018.
604 605	Heinrich Jiang and Ofir Nachum. Identifying and correcting label bias in machine learning. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 702–712. PMLR, 2020.
606 607 608 609	Sangwon Jung, Donggyu Lee, Taeeon Park, and Taesup Moon. Fair feature distillation for visual recognition. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 12115–12124, 2021.
610 611 612	Sangwon Jung, Sanghyuk Chun, and Taesup Moon. Learning fair classifiers with partially annotated group labels. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10348–10357, 2022.
613 614 615 616	Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoo Yun, Hwanjun Song, Joonhyun Jeong, Jung- Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameterization. In <i>International Conference on Machine Learning</i> , pp. 11102–11118. PMLR, 2022.
617 618 619	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 4015–4026, 2023.
620 621	Yann LeCun, Corinna Cortes, Chris Burges, et al. Mnist handwritten digit database, 2010.
622 623 624	Dong Bok Lee, Seanie Lee, Joonho Ko, Kenji Kawaguchi, Juho Lee, and Sung Ju Hwang. Self- supervised dataset distillation for transfer learning. In <i>The Twelfth International Conference on</i> <i>Learning Representations</i> , 2023.
625 626 627	Hae Beom Lee, Dong Bok Lee, and Sung Ju Hwang. Dataset condensation with latent space knowledge factorization and sharing. <i>arXiv preprint arXiv:2208.10494</i> , 2022a.
628 629 630	Saehyung Lee, Sanghyuk Chun, Sangwon Jung, Sangdoo Yun, and Sungroh Yoon. Dataset condensa- tion with contrastive signals. In <i>International Conference on Machine Learning</i> , pp. 12352–12364. PMLR, 2022b.
632 633 634	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre- training for unified vision-language understanding and generation. In <i>International conference on</i> <i>machine learning</i> , pp. 12888–12900. PMLR, 2022.
635 636 637 638	Evan Zheran Liu, Behzad Haghgoo, Annie S. Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. <i>ArXiv</i> , abs/2107.09044, 2021. URL https://api.semanticscholar.org/CorpusID:235825419.
640 641	Songhua Liu, Kai Wang, Xingyi Yang, Jingwen Ye, and Xinchao Wang. Dataset distillation via factorization. <i>arXiv preprint arXiv:2210.16774</i> , 2022.
642 643 644 645	Yanqing Liu, Jianyang Gu, Kai Wang, Zheng Hua Zhu, Wei Jiang, and Yang You. Dream: Efficient dataset distillation by representative matching. 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 17268–17278, 2023. URL https://api.semanticscholar.org/CorpusID:257232785.
647	Noel Loo, Ramin Hasani, Alexander Amini, and Daniela Rus. Efficient dataset distillation using random feature approximation. <i>arXiv preprint arXiv:2210.12067</i> , 2022.

648 649 650	Noel Loo, Ramin Hasani, Mathias Lechner, Alexander Amini, and Daniela Rus. Understanding reconstruction attacks with the neural tangent kernel and dataset distillation. <i>arXiv preprint arXiv:2302.01428</i> , 2023.
651 652 653	Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. <i>arXiv preprint arXiv:1511.00830</i> , 2015.
654 655 656	Yao Lu, Jianyang Gu, Xuguang Chen, Saeed Vahidian, and Qi Xuan. Exploring the impact of dataset bias on dataset distillation. <i>ArXiv</i> , abs/2403.16028, 2024. URL https://api. semanticscholar.org/CorpusID:268680434.
658 659	Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. <i>ACM computing surveys (CSUR)</i> , 54(6):1–35, 2021.
660 661 662	Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely wide convolutional networks. <i>Advances in Neural Information Processing Systems</i> , 34:5186–5198, 2021.
664 665 666	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. <i>arXiv preprint arXiv:2304.07193</i> , 2023.
667 668 669	Novi Quadrianto, Viktoriia Sharmanska, and Oliver Thomas. Discovering fair representations in the data domain. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 8227–8236, 2019.
670 671 672 673 674	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021.
675 676 677	Harsh Rangwani, Sumukh K Aithal, Mayank Mishra, and R. Venkatesh Babu. Escaping saddle points for effective generalization on class-imbalanced data. <i>ArXiv</i> , abs/2212.13827, 2022. URL https://api.semanticscholar.org/CorpusID:255186001.
678 679 680	Prasanna Sattigeri, Samuel C Hoffman, Vijil Chenthamarakshan, and Kush R Varshney. Fairness gan: Generating datasets with fairness properties using a generative adversarial network. <i>IBM Journal</i> <i>of Research and Development</i> , 63(4/5):3–1, 2019.
681 682 683 684 685	Shivashankar Subramanian, Afshin Rahimi, Timothy Baldwin, Trevor Cohn, and Lea Frermann. Fairness-aware class imbalanced learning. In <i>Conference on Empirical Methods in Natural</i> <i>Language Processing</i> , 2021. URL https://api.semanticscholar.org/CorpusID: 237593022.
686 687	Ilia Sucholutsky and Matthias Schonlau. Soft-label dataset distillation and text dataset distillation. In 2021 International Joint Conference on Neural Networks (IJCNN), pp. 1–8. IEEE, 2021.
688 689 690 691	Davoud Ataee Tarzanagh, Bojian Hou, Boning Tong, Qi Long, and Li Shen. Fairness-aware class imbalanced learning on multiple subgroups. <i>Proceedings of machine learning research</i> , 216:2123–2133, 2023. URL https://api.semanticscholar.org/CorpusID:260843665.
692 693 694	Saeed Vahidian, Mingyu Wang, Jianyang Gu, Vyacheslav Kungurtsev, Wei Jiang, and Yiran Chen. Group distributionally robust dataset distillation with risk minimization. <i>arXiv preprint arXiv:2402.04676</i> , 2024.
695 696 697 698	Robin Vogel, Mastane Achab, Stéphan Clémençon, and Charles Tillier. Weighted empirical risk mini- mization: Sample selection bias correction based on importance sampling. ArXiv, abs/2002.05145, 2020. URL https://api.semanticscholar.org/CorpusID:211082802.
699 700 701	Kai Wang, Bo Zhao, Xiangyu Peng, Zheng Zhu, Shuo Yang, Shuo Wang, Guan Huang, Hakan Bilen, Xinchao Wang, and Yang You. Cafe: Learning to condense dataset by aligning features. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12196–12205, 2022.

702 703 704 705	Shaobo Wang, Yantai Yang, Qilong Wang, Kaixin Li, Linfeng Zhang, and Junchi Yan. Not all samples should be utilized equally: Towards understanding and improving dataset distillation, 2024. URL https://arxiv.org/abs/2408.12483.
706 707	Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. <i>arXiv</i> preprint arXiv:1811.10959, 2018.
708 709 710 711	Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. Towards fairness in visual recognition: Effective strategies for bias mitigation. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 8919–8928, 2020.
712 713 714 715	Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rogriguez, and Krishna P Gummadi. Fairness constraints: Mechanisms for fair classification. In <i>Artificial intelligence and statistics</i> , pp. 962–970. PMLR, 2017.
716 717	Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In <i>International conference on machine learning</i> , pp. 325–333. PMLR, 2013.
718 719 720 721	Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In <i>Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society</i> , pp. 335–340, 2018.
722 723	Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In <i>International Conference on Machine Learning</i> , pp. 12674–12685. PMLR, 2021.
724 725	Bo Zhao and Hakan Bilen. Synthesizing informative training samples with gan. In <i>NeurIPS 2022</i> <i>Workshop on Synthetic Data for Empowering ML Research</i> , 2022.
727 728	Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pp. 6514–6523, 2023.
729 730 731	Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching. <i>arXiv preprint arXiv:2006.05929</i> , 2020.
732 733	Zhenghao Zhao, Haoxuan Wang, Yuzhang Shang, Kai Wang, and Yan Yan. Distilling long-tailed datasets, 2024. URL https://arxiv.org/abs/2408.14506.
734 735 736 727	Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. Anomalyclip: Object-agnostic prompt learning for zero-shot anomaly detection. In <i>The Twelfth International Conference on Learning Representations</i> , 2023.
738 739	Yongchao Zhou, Ehsan Nezhadarya, and Jimmy Ba. Dataset distillation using neural feature regression. <i>arXiv preprint arXiv:2206.00719</i> , 2022.
740	
741	
742	
743	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	

Datasets	TA	PA	TA number	PA number	Training set size	Test set size	BR in Training set	BR in Test set	Con 10	densed 1 50	ratio 100
C-MNIST (FG)	Digital number	Digital color	10	10	60000	10000	0.90	balance	0.17%	0.83%	1.67%
C-MNIST (BG)	Digital number	Background color	10	10	60000	10000	0.90	balance	0.17%	0.83%	1.67%
C-FMNIST (FG)	Object category	Object color	10	10	60000	10000	0.90	balance	0.17%	0.83%	1.67%
C-FMNIST (BG)	Object category	Background color	10	10	60000	10000	0.90	balance	0.17%	0.83%	1.67%
CIFAR10-S	Object category	Grayscale or not	10	2	50000	20000	0.90	balance	0.20%	1.00%	2.00%
CelebA	Attractive	Gender	2	2	162770	7656	class0: 0.62 class1: 0.77	balance	0.012%	0.061%	0.12%

#### Table 6: Statistics for all datasets used in our paper.

#### DATASEST STATISTICS А

767 In this section, we provide detailed statistics for all datasets used in the manuscript for reproduction. 768 As shown in Table 6, we present the target attribute (TA), protected attribute (PA), the sample number 769 of the training set, the sample number of the test set, and the BR in the training set. Additionally, all test sets are balanced, with equal sample sizes across groups. We also report the condensed ratio at 770 IPC 10, 50, and 100, which is computed by the ratio of the condensed dataset size to the training set 771 size. 772

#### В MORE ATTRIBUTES ANALYSIS ON CELEBA

We explore additional facial attributes in CelebA to further demonstrate the robustness of FairDD. To this end, we regard gender as the PA, and young, big\_nose, and blond\_hair as the TA, 778 which results in CelebA<sub> $\mu$ </sub>, CelebA<sub>h</sub>, CelebA<sub>h</sub> and respectively. We also replace the target attribute</sub></sub>from attribute to  $blond_hair$ , resulting in CelebA<sup>h</sup>. The performance is reported in Table 7. FairDD exhibits fairer behavior compared to vanilla DDs across these attributes while maintaining comparable accuracy, as seen in Table 8.

Tabl	Table 7: Fairness comparison on different attributes.												Table	8:	Ac	curacy	y cor	npari	son.
Methods Dataset	IPC	$\left  \frac{D}{DEO_M} \right $	M DEO <sub>A</sub>	DM+F	airDD DEO <sub>A</sub>	DEO <sub>M</sub>	C DEO <sub>A</sub>	DC+F	airDD DEO <sub>A</sub>	$\left \frac{W}{DEO_M}\right $	DEO <sub>A</sub>		Methods Dataset	IPC	$\left  \frac{\mathrm{DM}}{\mathrm{Acc.}} \right $	+FairDD Acc.	$\left  \frac{DC}{Acc.} \right $	+FairDD Acc.	Whole Acc.
$\operatorname{CelebA}_y$	10 50 100	34.18 46.90 44.96	31.49 41.13 37.84	13.30 12.90 9.17	10.38 8.21 5.11	20.58 27.98 27.76	19.26 25.18 24.26	10.86 14.69 19.03	8.55 11.26 13.61	25.40	16.02		$\operatorname{CelebA}_y$	10 50 100	62.34 63.59 66.68	63.79 67.33 69.90	55.91 59.87 63.53	<b>56.99</b> 59.42 61.59	75.99
CelebA <sub>b</sub>	10 50 100	45.57 51.91 52.75	45.13 51.13 51.27	15.63 14.44 8.03	13.47 12.01 6.10	18.17 23.85 24.48	16.81 22.34 23.53	7.54 20.58 12.15	6.34 16.87 11.00	34.48	25.50		CelebA <sub>b</sub>	10 50 100	57.46 58.71 60.30	59.50 62.39 64.34	52.91 56.55 57.65	<b>54.67</b> 55.46 57.15	66.80
$CelebA_h$	10	17.01	9.56	7.76	6.02	12.44	8.01	9.25	7.31	15.53	11.56		$CelebA_h$	10	63.64	64.86	58.04	57.55	75.33
$CelebA^h$	10	30.28	20.76	12.70	8.28	25.94	15.11	16.78	9.88	46.67	26.11		$CelebA^h$	10	77.66	79.71	72.07	75.03	79.44

To further study the generalization of FairDD, we regard blond hair as the protected attribute and attractive as the target attribute, resulting in  $CelebA_b$ . As illustrated in Table, FairDD+DM obtains 7.76% DEO<sub>M</sub> and 6.02% DEO<sub>A</sub>, outperforming DM by 9.25% and 3.54%. Accuracy has also been improved.

#### С **COMPUTATION OVERHEAD**

800 In this section, we investigate the computational efficiency of FairDD. Since FairDD performs finegrained alignment at the group level, we evaluate the impact of the number of groups on training 801 time (min) and peak GPU memory consumption (MB). As shown in Table 9, FairDD requires more 802 time than vanilla DDs on C-MNIST (FG), and the time increases as the number of groups (PA) 803 grows. This phenomenon is particularly noticeable in DC because the gradient must be computed 804 with respect to the number of groups. In contrast, DM avoids computing gradients independently 805 for each group and directly computes all embeddings once in a single pass, reducing the additional 806 overhead caused by FairDD's group-level alignment. Regarding GPU memory usage, FairDD incurs 807 no obvious additional overhead compared to vanilla DDs. 808

Here, we further supplement the overhead analysis with respect to image resolutions. We conduct 809 experiments on CMNIST, CelebA (32), CelebA (64), and CelebA (96) on DM and DC at IPC=10.

756 758

766

773 774

775 776

777

779

791 792

793

794

795

796 797 798

Table 9: Comparison of computation overhead on FairDD and vanilla DDs.

Group	0 (vani	lla DD)	2 (Fa	irDD)	4 (Fa	irDD)	6 (Fa	irDD)	8 (Fa	irDD)	10 (Fa	airDD)
number	T (min)	$G\left(MB ight)$	T (min)	$G\left(MB\right)$	T (min)	G (MB)	T (min)	$G\left(MB\right)$	T (min)	G (MB)	T (min)	G (MB)
DC	70	2143	94	2345	128	2369	152	2393	181.8	2419	210	2443
DM	26.2	1579	31.75	1579	33.2	1579	35.2	1579	36.5	1579	36.9	1579

DM and DC align different signals, which would bring different effects. As illustrated in Table 10, it can be observed that FairDD + DM does not require additional GPU memory consumption but does necessitate more time. The time gap increases from 0.42 minutes to 1.79 minutes as input resolution varies (e.g., CelebA  $32 \times 32$ , CelebA  $64 \times 64$ , and CelebA  $96 \times 96$ ); however, the gap remains small. This can be attributed to FairDD performing group-level alignment on features, which is less influenced by input resolution. Notably, although CMNIST and CelebA  $(32 \times 32)$  share the same resolution, the time gap is more pronounced for CMNIST (e.g., 3 minutes). This is attributed to CMNIST having 10 attributes, whereas CelebA  $(32 \times 32)$  has only 2 attributes. These indicate that FairDD + DM requires no additional GPU memory consumption. Its additional time depends on both input resolution and the number of groups, but the number of groups more significantly influences it. As for DC, FairDD requires additional GPU memory and time. Since FairDD + DC explicitly computes group-level gradients, the resulting gradient caches cause FairDD + DC to consume more memory. The additional consumption is acceptable compared to the performance gain on fairness. Additionally, the time gap is relatively larger than that observed between DM and FairDD + DM. Similar to DM, the group number is the primary factor contributing to additional time consumption compared to input resolution. 

Table 10: Comparison of computation overhead for IPC = 10.

Methods	Group	1	DM	DM+	-FairDD		DC	DC+	FairDD
Dataset	number	Time	Memory	Time	Memory	Time	Memory	Time	Memory
CMNIST	10	15.55	1227	18.55	1227	58.75	1767	83.13	1893
$CelebA32 \times 32$	2	10.93	2293	11.35	2293	32.98	2413	34.65	2479
$CelebA64 \times 64$	2	11.18	8179	12.20	8177	43.67	8525	47.07	8841
$CelebA96 \times 96$	2	12.83	17975	14.62	17975	82.37	18855	86.88	19437

## D ADDITIONAL EXPERIMENT ON UTKFACE

We also conduct another real image dataset namely UTKFace, commonly used for fairness. It consists of 20,000 face images including three attributes, age, gender, and race. We follow a common setting and treat age as the target attribute and gender as the protected attribute. We test DM and FairDD + DM with the same parameters, the results in Table 11 show that our method outperforms the vanilla dataset distillation by 16.1% and 8.92% on the DEO<sub>M</sub> and DEO<sub>A</sub>. Similar results are observed at IPC = 50.

Table 11: UTKFace performance on DM.

Methods Dataset	IPC	Acc.	DM DEO <sub>M</sub>	DEOA	$\frac{D}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>
UTKFace	10	66.67	32.97	16.31	67.72	16.87	7.39
	50	73.15	28.58	14.03	74.66	10.59	5.09

#### E ABLATION STUDY ON WEIGHTING MECHANISM

860 We denote the model with inversely proportional weighting as  $\text{FairDD}_{\text{inverse}}$ . Our experiments on 861 C-FMNIST and CIFAR10-S at IPC=10 reveal that  $\text{FairDD}_{\text{inverse}}$  suffers significant performance 862 degradation, with  $\text{DEO}_M$  increasing from 33.05% to 56.60% and  $\text{DEO}_A$  rising from 19.72% to 863 35.13% in terms of fairness performance metrics. Additionally, there is a decline in accuracy for TA.

Methods Dataset	IPC	Acc.	$\frac{DM}{DEO_M}$	DEOA	$\frac{\rm DM}{\rm Acc.}$	FairDI	Dinverse DEO <sub>A</sub>	$\frac{DN}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>
C-FMNIST (BG)	10	22.26	100.0	99.05	70.55	56.60	35.13	71.10	33.05	19.72
CIFAR10-S	10	37.88	59.20	39.31	38.14	48.27	37.41	45.17	31.75	8.73

Table 12: Performance between FairDD and FairDD<sub>inverse</sub>.

We attribute this degradation to the excessive penalization of groups with larger sample sizes. The success of FairDD lies in grouping all samples with the same PA into a single group and performing the group-level alignment. Each group contributes equally to the total alignment, inherently mitigating the effects of imbalanced sample sizes across different groups. However, penalizing groups based on sample cardinality reintroduces an unexpected bias related to group size in the information condensation process. This results in large groups receiving smaller weights during alignment, placing them in a weaker position and causing synthetic samples to deviate excessively from large (majority) groups. Consequently, majority patterns become underrepresented, ultimately hindering overall performance.

### F ABLATION STUDY ON GROUP LABEL NOISE

Here, we evaluate the robustness of spurious group labels could provide more insights. We randomly sample the entire dataset according to a predefined ratio. These samples are randomly assigned to group labels to simulate noise. To ensure a thorough evaluation, we set sample ratios at 10%, 15%, 20%, and 50%. As shown in the table, when the ratio increases from 10% to 20%, the DEO<sub>M</sub> results range from 14.93% to 18.31% with no significant performance variations observed. These results indicate that FairDD is robust to noisy group labels. However, as the ratio increases further to 50%, relatively significant performance variations become apparent. It can be understood that under a high noise ratio, the excessive true samples of majority attributes are assigned to minority labels. This causes the minority group center to shift far from its true center and thus be underrepresented.

#### Table 13: Ablation study on group label noise.

Methods Dataset	IPC Acc	DM	DEOA	$\frac{D}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>	$\frac{\text{DM}+}{\text{Acc.}}$	FairDD DEO <sub>M</sub>	$\frac{(10\%)}{\text{DEO}_{A}}$	$\frac{\text{DM}+}{\text{Acc.}}$	FairDD DEO <sub>M</sub>	(15%) DEO <sub>A</sub>	$\frac{DM+}{Acc.}$	FairDD DEO <sub>M</sub>	(20%) DEO <sub>A</sub>	$\frac{\text{DM}+}{\text{Acc.}}$	FairDD DEO <sub>M</sub>	(50%) DEO <sub>A</sub>
CMNIST (BC	G)  10  27.9	5 100.0	99.11	94.88	13.42	6.77	94.34	16.54	7.81	94.44	17.90	8.61	94.32	18.31	9.20	89.56	66.19	25.97

#### G ABLATION STUDY ON BALANCED ORIGINAL DATASET.

We synthesized a fair version of CelebA, referred to as CelebA<sub>*Fair*</sub>. The target attribute is attractive (attractive and unattractive), and the protected attribute is gender (female and male). In the original dataset, the sample numbers for female-attractive, female-unattractive, male-attractive, and maleunattractive groups are imbalanced. To create a fair version, Celeb $A_{Fair}$  samples the number of instances based on the smallest group, ensuring equal representation across all four groups. We tested the fairness performance of FairDD and DM at IPC = 10, as well as the performance of models trained on the full dataset. As shown in Table 14, vanilla DD achieves 14.33% DEO<sub>A</sub> and 8.77% DEO<sub>M</sub>. In comparison, the full dataset achieves 3.66% DEO<sub>A</sub> and 2.77% DEO<sub>M</sub>. While DM still exacerbates bias with a relatively small margin, this is primarily due to partial information loss introduced during the distillation process. FairDD produces fairer results, achieving 11.11% DEO<sub>A</sub> and 6.68% DEO<sub>M</sub>. 

Methods Dataset	IPC	Acc.	Whole DEO <sub>M</sub>	DEO <sub>A</sub>	Acc.	DM DEO <sub>M</sub>	DEOA	$\frac{D}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>
CelebA <sub>Fair</sub>	10	76.33	3.66	2.77	63.31	14.33	8.77	63.17	11.11	6.68

Table 14: Performance on balanced original dataset.

### H ABLATION STUDY ON NUANCED PA GROUPS

We perform a fine-grained PA division. For example, we consider gender and wearing-necktie as two correlated attributes and divide them into four groups: males with a necktie, males without a necktie, females with a necktie, and females without a necktie (CelebA<sub>g&n</sub>). Similarly, we consider gender and paleskin and divide them into four groups (CelebA<sub>g&p</sub>). Their target attribute is attractive. As shown in the Table 15, FairDD outperforms vanilla DD in the accuracy and fairness performance on these two experiments. The performance for necktie and gender is improved from 57.50% to 25.00% on DEO<sub>M</sub> and 52.79% to 21.73% on DEO<sub>A</sub>. Accuracy is also improved from 63.25% to 67.98%. Similar results can be observed for gender and paleskin. Hence, FairDD can mitigate more fine-grained attribute bias, even when there is an intersection between attributes.

Table 15: Performance on nuanced groups.

Methods Dataset	IPC	Acc.	DM DEO <sub>M</sub>	DEOA	$\frac{D}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>
${f CelebA}_{g\&n}$	10	63.25	57.50	52.79	67.98	25.00	21.73
${f CelebA}_{g\&p}$	10	62.48	44.81	41.60	64.37	26.92	19.33

## I ABLATION STUDY ON IMBALANCED PA GROUPS

To further study FairDD robustness under more biased scenarios, we keep the sample number of the majority group in each class invariant and allocate the sample size to the remaining 9 minority groups with increasing ratios, i.e., 1:2:3:4:5:6:7:8:9. We denote this variant CMNIST<sub>unbalance</sub> This could help create varying extents of underrepresented samples for different minority groups. Notably, the least-represented PA groups account for only about 1/500 of the entire dataset, which equates to just 12 samples out of 6000 in CMNIST<sub>unbalance</sub>. As shown in Table 16, FairDD achieves a robust performance of 16.33% DEO<sub>M</sub> and 9.01% DEO<sub>A</sub> compared to 17.04% and 7.95% in the balanced PA groups. A similar steady behavior is observed in accuracy, which changes from 94.45% to 94.61%. This illustrates the robustness of FairDD under different levels of dataset imbalance.

Table 16: Performance on imbalanced PA.

Methods Dataset	IPC	Acc.	DM DEO <sub>M</sub>	DEO <sub>A</sub>	$\frac{D}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>
CMNIST	10	25.01	100.0	99.96	94.61	17.04	7.95
CMNIST <sub>unbalance</sub>	10	23.38	100.0	99.89	94.45	16.33	9.01

962 963

918

925 926

927

928

929

930

931

932

933

934 935

945 946

947

948

949

950

951

952

953

954 955

#### 964 965

966

## J EXPLORATION ON VISION TRANSFORMER AS BACKBONE

Although the Vision Transformer (ViT) is a powerful backbone network, to the best of my knowledge, current DDs, such as DM and DC, have not yet utilized ViT as the extraction network. We conducted experiments using 1-layer, 2-layer, and 3-layer ViTs. As shown in Table 17, vanilla DM at IPC=10
suffers performance degradation in classification, dropping from 25.01% to 18.63%. Moreover, as the number of layers increases, the performance deteriorates more severely. This suggests that current DDs are not directly compatible with ViTs. While FairDD still outperforms DM in both accuracy

V

and fairness metrics, the observed improvement gain is smaller compared to results obtained on convolutional networks. Further research into leveraging ViTs for DD and FairDD is a promising direction worth exploring. 

	Table 17:	Exploration	on ViT	architecture.
--	-----------	-------------	--------	---------------

Methods Dataset	IPC	Acc.	DM DEO <sub>M</sub>	DEOA	$\frac{DI}{Acc.}$	M+Fair DEOм	$\frac{DD}{DEO_A}$
ViT1	10	18.63	100.0	98.48	56.15	82.10	56.72
ViT2 ViT3	10 10	18.28 16.15	100.0 100.0	98.99 95.75	33.89 26.70	72.85 65.71	40.97 29.46

#### Κ **EXPLORATION ON CHALLENGING DATASET**

We created CIFAR100-S following the same operation as CIFAR10-S, where the grayscale or not is regarded as PA. Due to the time limit, we supplemented CIFAR100-S on DM at IPC=10. As shown in Table 18, DM achieves the classification accuracy of 22.84%, and the fairness of 69.9%  $DEO_M$  and 25.37% DEO<sub>A</sub>. Compared to vanilla DM, FairDD obtains more accurate classification performance and mitigates the bias to the minority groups, with 27.30% DEO<sub>M</sub> and 15.54% DEO<sub>A</sub> improvement.

Table 18: Exploration on challenging dataset.

Methods Dataset	IPC	Acc.	DM DEO <sub>M</sub>	DEOA	$\frac{DN}{Acc.}$	M+Fair DEO <sub>M</sub>	DD DEO <sub>A</sub>
CIFAR100-S	10	19.69	69.90	25.37	22.84	42.60	9.83

#### L COMPARISON WITH MTT

Unlike DMF, MTT uses a two-stage method to condense the dataset. First, it stores the model trajectories, and then it uses these trajectories to guide the generation of the synthetic dataset. To provide a comprehensive comparison, we compare FairDD with MTT, as shown in Tables 19 and 20.

Table 19: Fairness comparison on diverse IPCs. The best Table 20: Accuracy comparison on diresults are high 

100.0 96.05 97.20 63.66 100.0 93.88 **21.95 13.33** 

e highlighte	ed in bold.		verse IPCs.	
$\operatorname{IPC} \left  \frac{\operatorname{Ran}\operatorname{dom}}{\operatorname{DEO}_{\mathrm{M}}} \frac{\operatorname{DEO}_{\mathrm{A}}}{\operatorname{DEO}_{\mathrm{A}}} \right $	$\left  \frac{MTT}{DEO_M} \right  \frac{DM}{DEO_A} \left  \frac{DM}{DEO_M} \right $	$\frac{DM+FairDD}{DEO_A} \frac{DM+FairDD}{DEO_A} \frac{W}{DEO_A}$	$\frac{\text{Whole}}{\text{DEO}_{A}} \qquad \frac{\text{Methods}}{\text{Dataset}} \left  \text{IPC} \right  \frac{\text{Rando}}{\text{Acc.}}$	$ \begin{array}{c c} m & MTT \\ \hline Acc. \end{array} & \hline DM \\ \hline Acc. \end{array} & + FairDD \\ \hline Acc. \end{array} & \hline Whole \\ \hline Acc. \end{array} $
10   100.0 98.72 50   100.0 99.58 100   100.0 88.64	25.70         14.86         100.0         9           25.46         12.60         100.0         9           26.81         13.02         99.36         6	99.96     17.04     7.95       01.68     10.05     5.46       16.38     8.17     4.86	0 5.89 C-MNIST   10   30.75 (FG)   10   30.75 50   47.38 100   67.41	92.00         25.01         94.61           94.08         56.84         96.58         97.71           94.29         78.04         96.79         97.71
10 100.0 99.40 50 100.0 98.52	97.00 62.46 100.0 9 96.60 62.02 100.0 9	99.68 <b>33.05 19.72</b> 99.71 <b>24.50 14.47</b> 91.40	0 51.68 C-FMNIST 10 24.96 (BG) 50 34.92	6 67.92 22.26 <b>71.10</b> 2 70.32 36.27 <b>79.07</b> 77.97

100 44.87 70.74 49.30 80.63

Methods

Dataset

C-MNIST

(FG)

C-FMNIST

(BG)

IPC DE

From the results, FairDD outperforms MTT in both fairness and accuracy. Notably, MTT surpasses DM by a large margin, which we attribute to two factors: 1) Unlike DMF, which is directly influenced by biased data, MTT aligns the model parameters to optimize the synthetic dataset, and this indirect alignment reduces the impact of bias in the data. 2) An accurate model typically conceals its inherent unfairness, as it can better classify each class despite underlying biases. For example, when Whole model achieves high accuracy on the C-MNIST (FG) dataset, MTT inherits this accuracy and conceals its biases. However, when the model's accuracy declines on the C-FMNIST (BG) dataset, MTT reveals its underlying unfairness in Fig. 6(a). In contrast, FairDD directly addresses unfairness rather than relying on high accuracy to obscure biased behavior in Fig. 6(b).



Figure 6: Visualization comparison on C-FMNIST (BG) between MTT and FairDD + DM.



Figure 7: Visualization on S at IPC=1 for FairDD and vanilla DDs. Left is the condensed dataset using FairDD, which incorporates different PA, i.e., foreground colors. **Right** is the condensed dataset using vanilla DDs, where each class presents the same color as the corresponding majority group.

#### ADDITIONAL VISUALIZATION ANALYSIS Μ

**Visualization analysis on**  $\mathcal{S}$  generation We aim to investigate whether FaiDD renders the expecta-tion of S locate the center among all groups, as clarified in Eq. 8. If the clarification holds, S should contain all PA at IPC = 1 because the expectation of S is equal to S when IPC =1. We visualize S at IPC=1 on C-MNIST (FG), where each class (digital number) is dominated by one color, and the rest is colored by the rest nine colors. As shown in Fig. 7, the S generated by FairDD combines all colors from PA groups. This suggests that FairDD can effectively incorporate all PA into resulting S, indirectly validating the Theorem 4.1. Meanwhile, we observe that the majority groups dominate vanilla DDs according to Eq. 6, where the resulting  $\mathcal{S}$  contains the colors from the corresponding majority groups.

#### Ν MORE VISUALIZATIONS

We provide more visualizations at IPC = 50 on different datasets in Figures 8, 9, 10, 11, and 13.













Figure 13: FVisualization comparison on CelebA between vanilla DM and FairDD + DM.

## $^{1404}_{1405}$ O Complete results in the format of mean $\pm$ standard deviation.

4.400			
1406			
1407			
1408			
1409			
1410			
1411			
1/12			
1/12			
1413			
1414			
1415			
1416			
1417			
1418			
1419			
1420			
1421			
1492			
1/22			
1423			
1424			
1425			
1426			
1427			
1428			
1429			
1430			
1431			
1432			
1433			
1434			
1/135			
1/26			
1407			
1437			
1438			
1439			
1440			
1441			
1442			
1443			
1444			
1445			
1446			
1447			
1448			
1449			
1450			
1450			
1451			
1452			
1453			
1454			
1455			
1456			
1457			

1458						~	
1459	0 <sup>A</sup>	:0.29	:0.11	±0.95	±1.34	±0.32	±0.54
1460	le DE(	£.89	5.78 <sub>±</sub>	1.72	1.68	3.17 <sub>4</sub>	4.16
1461	Vho	40	74	75 4	85.5	84 3	90 2
1462	EO <sub>1</sub>	0±0	0 <sup>±0.</sup>	070	0±1	$^{2\pm0}$	2 <sup>±0</sup>
1463		10.1	6.7	79.2	91.4	49.7	24.8
1464	A	0.39 0.13 0.34	0.63 0.21 0.25	-0.87 1.20 1.00	0.64 1.02 0.94	0.35 0.80 0.80	0.89 0.88 1.26
1465	DE(	.21± .50± .88±	.83± .04± .80±	0.87 <sub>±</sub> 33± 40±	1.50± 1.70± 1.20±	.77± .53± .06±	.36± +29±
1466	airl 1	797 537 246	56 6 22 9 12 7	45 2( 00 22 22 40	48 23 37 23 48 23	26 7 29 6 71 6	20 5 54 14 82 15
1467	H OE	$^{8\pm 0.3}_{7\pm 1.1}$	$1\pm 0.1$ $1\pm 1.1$ $3\pm 1.1$	$^{0\pm 2.5}_{7\pm 0.6}$	0 <sub>±2.</sub> , 0 <sub>±3.</sub> ;	$^{0\pm 2.5}_{\pm 1.1}$	)±1.5 2±1.5 7±1.3
1468	D	11.8 18.3 11.8	15.3 20.4 16.3	42.00 47.6'	36.0 34.0	20.80 14.70 12.10	9.2( 23.6) 24.8'
1469 <sup>p</sup> p	$O_{\rm A}$	E2.96 E2.52 E1.55	E2.16 E2.94 E2.22	E1.20 E1.87 E2.25	E0.76 E1.30 E1.56	E0.39 E0.59 E0.91	E1.31 E0.69 E0.50
1470 Lug	AM DE	8.71 <sub>=</sub> 6.63 <sub>=</sub> 3.30 <sub>=</sub>	2.06 <sub>-1</sub> 3.30 <sub>-1</sub> 2.86 <sub>-1</sub>	5.17 <sub>1</sub> 3.27 <sub>1</sub> 0.10 <sub>1</sub>	4.70 8.40 3.90	6.19 <sub>4</sub> 6.19 <sub>4</sub> 8.12 <sub>4</sub>	6.70 <sub>1</sub> 8.53 <sub>1</sub> 9.34
1471	ME	96 7 80 2 87 3	2.50 5 .50 2 34 2	00 9.00 9.00 9.00 9.00 9.00 9.00 9.00 9	00.00 8 0.00 8 0.00 8	8.18 <sup>3</sup> 31 <sup>3</sup> 16 <sup>4</sup>	353 603 843
1472	D	$99 \pm 1$ $03 \pm 1$ $30 \pm 1$	).0±2 15±1 88±1	).0±0 67±0 33±0	0.0±0 0.0±0.0	$80 \pm 3$ $80 \pm 1$ $30 \pm 2$	75±1 57±1 51±0
1473 🗧		4 98. 2 52. 0 69.	<sup>6</sup> 100 1 64. 3 44.	5   100 7   99. 9   97.	8 100 4 100 2 100	7  51. 4  56. 1  82.	5 40. 4 43. 1 52.
1474 ee	DO <sub>A</sub> O	±0.4 ±0.4 ±0.3	±0.2 ±0.3	上1.4 土0.8 土0.7	±0.2 ±1.5 ±2.6	±0.8' ±1.0' ±1.6	±1.7 ±0.9 ±1.6
1475 <sup>1</sup> <sub>JC</sub>	DE DE	6.64 5.94 6.21	7.50. 8.05. 6.32.	16.80 17.37 17.43	23.80 34.50 36.00	5.28. 9.10. 25.18	4.49 15.15 17.63
1476	± t ≱	L.18 26 L.65	1.52 2.48 1.38	L.22 2.15 L.37	1.59 2.87 2.58	2.95 L.75 2.97	2.37 1.37 1.65
1477	DEO DEO	24±1 18±1 88±1	$61\pm 1$ $66\pm 2$ $24\pm 1$	$33\pm 1$ $00\pm 2$ $33\pm 1$	1 <sup>4</sup> 00	5±06 1±00 2±08	36±2 57±1 27±1
1478 <sup>U</sup>		10 12 40 9. 59 11	47 18 87 19 42 12	44 32 81 40 04 48	25 42 43 78 82 77	20 19 78 29 51 62	45 8. 43 22 34 28
1479 4	EO,	$5\pm 1.$ $1\pm 1.$ $2\pm 1.1$	$0\pm 2$ , $3\pm 1$ , $9\pm 1$ ,	7±1. 3±0.3 0±2.0	$0\pm 0.1$ $0\pm 1.4$ $0\pm 1.4$	$8 \pm 1.5$ $6 \pm 0.5$ $3 \pm 0.1$	$8 \pm 1.5$ $0 \pm 1.5$ $6 \pm 0.3$
1480	DIG	) 91.4 3 34.9 3 35.8	) 88.3 5 42.2 7 27.7	97.2 () 81.9 () 79.1	) 99.4 ) 95.6 ) 90.7	3 48.3   60.5 2 60.9	1 34.4 5 46.5 2 42.6
1481 9		±0.00 ±2.98 ±1.68	±0.00 ±3.35 ±1.37	±0.00 ±0.00	±0.00 ±0.00	±2.85 ±0.61 ±0.92	±1.51 ±2.05 ±1.05
1482 Od	DE	100.0_ 55.34_ 54.36_	100.0 33.05 53.24	0.001	100.0	30.70 92.00	34.85 56.74 50.99
1483 <sup>1</sup>		.59 <sup>]</sup> .35 6	.57   54   64	.06 84 1 92	1 01. 1 09 1 00 00 00 00 00 00 00 00 00 00 00 00 0	9 00. 84 9 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	.62 .30 .93
1484	DD	96±0 86±0 <sup>19</sup> ±0	$ \frac{34\pm0}{90\pm0} $	$01_{\pm 1}$ $48_{\pm 0}$ $07_{\pm 1}$	$88 \pm 1$ $25 \pm 0$ $48 \pm 1$	$22 \pm 1$ $21 \pm 0$ $26 \pm 0$	29±1 92±1 94±2
1485 H	Fair	4 11.4 9 8.8 5 9.4	8 9.9 0 9.9 7 10.4	8 24. 3 21.	3 34.1 5 25.1 9 23.4	5 8.2 2 11.2 8 11.2	0 5.2 5 12.9 0 17.4
1486 Sal	OM OM	±2.1 ±1.3 ±1.8	±1.5 ±2.0	土2.95 土2.45 土1.90	±2.2 ±1.2	$\pm 1.1$ $\pm 2.2$ $\pm 1.5$	±1.60 ±1.31 ±2.81
1487	DE	26.75 18.42 22.32	20.66 20.29 32.58	46.80 46.67 56.68	51.75 44.60 52.75	22.08 34.39 32.70	6.64 14.33 18.16
1488	$\overline{A}$	2.78 1.03 0.99	2.31 1.59 1.28	1.98	0.95 1.20	1.52 2.23 1.17	2.59 3.18 2.21
1489	DEC	:.61± 1.55± 1.45±	.60± .38± .87±	1.96± 1.02± 1.45±	91± 41± 38±	'.35± .81± .64±	1.16± 1.83± 1.52±
1490	DO	00 65 39 20 53 17	00 73 05 26 85 20	47 78 50 67 37 66	00 92 00 75 37 73	83 27 00 45 53 48	62 14 32 23 42 27
1491	EO	$5\pm 0.$ $9\pm 1.$ $7\pm 1.$	$0\pm 0.$ $6\pm 3.$ $3\pm 1.$	$0\pm 0.$ $3\pm 0.$ $8\pm 0.$	0±0. 0±0.	$3\pm 1.$ $6\pm 1.$ $9\pm 1.$	8±2. 9±3. 0±2.
1492 App		99.8 46.9 45.2	100. 60.6 62.6	99.4 99.3	100.	42.2 71.4 68.6	15.4 24.8 29.0
1493	DA	=0.57 =0.40	=0.41 =0.22	±0.95 =0.68	=1.18 =0.66	=0.71 =0.59	=0.85 =1.00
1494 00	DE	7.95 <sub>4</sub> 5.46 <sub>4</sub> 4.86 <sub>4</sub>	6.77 <sub>±</sub> 5.25 <sub>±</sub> 4.31 <sub>±</sub>	16.38 <sub>-</sub> 3.74 <sub>4</sub> 2.75 <sub>4</sub>	9.72 <sub>4</sub> 4.47 <sub>4</sub> 3.33 <sub>4</sub>	8.73 <sub>4</sub> 7.35 <sub>4</sub> 5.89 <sub>4</sub>	5.71 <sub>±</sub> 9.87 <sub>±</sub> 6.65 <sub>±</sub>
1495	4 √	13 75 45	61 74 45	69 1 89 1 92 1	35 1 19 1 11 1	52 78 50	82 32 46
1490 3du	EO	<sup>4</sup> ±1. 5±0. 7±0.	$\frac{2 \pm 1}{8 \pm 0}$ .	$7\pm 1.$ $2\pm 1.$ $33\pm 1.$	$5\pm 2.$ $0\pm 1.$ $5\pm 1.$	5±1. 8±0. 7±0.	7±0. 8±0. 3±0.
1497	D	17.0 10.0	13.4 8.9 6.6	26.8 24.9 23.8	33.0 24.5 21.9	31.7 18.2 14.7	9.3 14.0
1490	$O_{\rm A}$	E0.05 E1.67 E1.80	E0.04 E1.11 E2.46	E0.14 E0.33 E1.09	E0.05 E0.10 E0.66	E1.28 E1.11 E1.45	E1.33 E0.70 E0.96
1500	$\frac{1}{DE}$	9.96 <sub>-</sub> 1.68 <sub>-</sub> 6.38 <sub>-</sub>	9.97 <sub>1</sub> 7.85 <sub>1</sub> 2.23	9.05 <sub>4</sub> 6.46 <sub>4</sub> 5.11 <sub>4</sub>	9.68 <sub>4</sub> 9.71 <sub>4</sub>	9.31 <sub>-</sub> 5.70 <sub>-</sub> 5.10 <sub>-</sub>	8.85 <sub>4</sub> 8.81 <sub>4</sub> 8.81 <sub>4</sub>
1501 H	M D	,00 9 ,00 9 ,27 6	9 00.1 9 00.1 2 00.1	9 00.1 9 00.1 8 00.1	9 00.1 9 00.1 9 00.1	27 <sup>3</sup> 24 <sup>5</sup> 11 <sup>5</sup>	342 943 184
1502	EO	).0±0 ).0±0 36±0	0±0.0 0±0.0	).0±0 0±0.0	0.0±0.0 0±0.0000.0000000000000000000000	20±1 13±1 81±2	01±1 26±0 63±1
1502 I		0 100 2 100 7 99.	100 100 100	1 100 8 100 4 100	4 100 3 100 7 100	5 59. 3 75. 5 73.	0 30. 2 40. 6 42.
1504 qr	<sup>1</sup> 30A	主 (1) (1) (1) (1) (1) (1) (1) (1) (1) (1)	±0.2 ±0.0 ±2.3	主0.1 主1.0 主0.7	土0.1 土0.6 土0.6	±0.6 '±1.8 '±1.4	±2.0 :±1.0 ±1.1
1505	DE	98.72 99.58 88.64	99.11 77.99 89.07	99.18 94.61 94.85	99.40 98.52 96.05	8.29 28.89 43.16	9.20 20.32 18.01
1506	$\frac{Ran}{M}$	0.00 0.00	0.00	0.00	0.00	0.93 1.42 1.88	1.97 1.41 1.19
1507	DEC	±0.00 ±0.00	±0.00 ±0.00	±0.00 ±0.00	±0.00 ±0.00	5.04± 7.11± 5.49±	).48± 2.88± 3.67±
1508		0 11 0 11 00 110	0 11 0 11 00 110	0 1c 0 1c	0 11 0 11 00 110	0  25 0  57 )0  66	0 1( 0 22 00 18
1509		$\frac{1}{5}$	$\frac{1}{5}$	$\frac{T_{5}}{100}$	$\frac{T_{5}}{100}$	S 5 11(	$\frac{1}{5}$
1510	nods tset	SIV (5	UIS' (E	SIN (E	SIN 3)	210-	þÅ
1511	Aetł: Data	(FC	(B(	FM. (FC	FM. (B(	FAF	Cele
	4 -	Ċ	Ċ	ڻ ان	Ċ	CI	

Table 22: Accuracy comparison on diverse IPCs. The results are reported in the format of mean $\pm$
standard deviation.

Methods Datasets	IPC	$\frac{\text{Random}}{\text{Acc.}}$	$\frac{\mathrm{DM}}{\mathrm{Acc.}}$	+FairDD Acc.	$\frac{DC}{Acc.}$	+FairDD Acc.	$\frac{\text{IDC}}{\text{Acc.}}$	+FairDD Acc.	DREAM Acc.	+FairDD Acc.	$\frac{\text{Whol}}{\text{Acc.}}$
C-MNIST (FG)	10 50 100	$\begin{array}{c} 30.75 {\pm} 0.96 \\ 47.38 {\pm} 0.98 \\ 67.41 {\pm} 1.08 \end{array}$	$25.01_{\pm 0.94} \\ 56.84_{\pm 1.92} \\ 78.04_{\pm 1.24}$	$94.61_{\pm 0.21} \\ 96.58_{\pm 0.08} \\ 96.79_{\pm 0.14}$	$\begin{array}{c} 71.41_{\pm 1.27} \\ 90.54_{\pm 0.43} \\ 91.64_{\pm 0.33} \end{array}$	$90.62_{\pm 0.28} \\ 92.68_{\pm 0.18} \\ 93.23_{\pm 0.19}$	$53.06_{\pm 1.13}$ $88.55_{\pm 0.38}$ $90.39_{\pm 0.48}$	$\begin{array}{c} 95.67_{\pm 0.20} \\ 96.77_{\pm 0.07} \\ 97.11_{\pm 0.08} \end{array}$	$75.04_{\pm 1.86} \\91.02_{\pm 0.68} \\88.87_{\pm 1.04}$	$\begin{array}{c} 94.04_{\pm 0.14} \\ 94.59_{\pm 0.23} \\ 95.16_{\pm 0.12} \end{array}$	97.71±0
C-MNIST (BG)	10 50 100	$27.95_{\pm 0.75} \\ 45.52_{\pm 0.98} \\ 67.28_{\pm 1.31}$	$23.40_{\pm 0.57} \\ 47.74_{\pm 1.35} \\ 79.87_{\pm 0.77}$	$94.88_{\pm 0.13} \\ 96.86_{\pm 0.09} \\ 97.33_{\pm 0.09}$	$\begin{array}{c} 65.91_{\pm 1.91} \\ 88.53_{\pm 0.61} \\ 90.20_{\pm 0.57} \end{array}$	$90.84_{\pm 0.19} \\92.20_{\pm 0.19} \\92.73_{\pm 0.17}$	$\begin{array}{c} 62.09 \\ \pm 1.06 \\ 86.14 \\ \pm 0.74 \\ 89.66 \\ \pm 0.66 \end{array}$	$\begin{array}{c} 94.84_{\pm 0.21} \\ 95.29_{\pm 0.13} \\ 95.84_{\pm 0.06} \end{array}$	$\begin{array}{c} 79.81_{\pm 1.37} \\ 89.24_{\pm 0.82} \\ 90.70_{\pm 0.70} \end{array}$	$\begin{array}{c} 93.54_{\pm 0.26} \\ 93.20_{\pm 0.50} \\ 94.06_{\pm 0.24} \end{array}$	97.80±
C-FMNIST (FG)	10 50 100	$\begin{array}{c} 32.80_{\pm 1.44} \\ 42.48_{\pm 1.05} \\ 55.31_{\pm 0.67} \end{array}$	$33.35_{\pm 1.27}$ $49.94_{\pm 0.75}$ $57.99_{\pm 0.84}$	$77.09 \pm 0.33$ $82.11 \pm 0.20$ $83.25 \pm 0.19$	$\begin{array}{c} 60.77 \pm 0.88 \\ 69.08 \pm 0.54 \\ 68.84 \pm 0.61 \end{array}$	$76.01_{\pm 0.19}$ $75.83_{\pm 0.33}$ $74.91_{\pm 0.40}$	${}^{44.08 \pm 1.64}_{64.45 \pm 0.69}_{66.37 \pm 0.26}$	$79.66_{\pm 0.21}$ $80.80_{\pm 0.33}$ $80.28_{\pm 0.19}$	$49.72_{\pm 1.07} \\ 65.69_{\pm 1.04} \\ 68.25_{\pm 0.94}$	$77.24_{\pm 0.15}$ $78.79_{\pm 1.07}$ $78.51_{\pm 0.55}$	82.94 <sub>±</sub>
C-FMNIST (BG)	10 50 100	$24.96 \pm 0.63 \\ 34.92 \pm 0.57 \\ 44.87 \pm 0.65$	$22.26_{\pm 0.77} \\ 36.27_{\pm 0.72} \\ 49.30_{\pm 0.58}$	$71.10{\pm}0.43$ $79.07{\pm}0.27$ $80.63{\pm}0.16$	$\begin{array}{c} 47.32 \pm 0.96 \\ 60.58 \pm 0.72 \\ 62.70 \pm 0.66 \end{array}$	$\begin{array}{c} 68.51 {\pm} 0.38 \\ 75.80 {\pm} 0.28 \\ 71.76 {\pm} 0.28 \end{array}$	$37.59 \pm 0.76$ $46.20 \pm 1.10$ $48.61 \pm 1.00$	$72.67_{\pm 0.20}$ $73.72_{\pm 0.43}$ $73.18_{\pm 0.57}$	$\begin{array}{c} 45.30 {\pm} 0.78 \\ 53.62 {\pm} 0.50 \\ 53.32 {\pm} 0.78 \end{array}$	$71.56_{\pm 0.34}$ $72.80_{\pm 0.11}$ $73.00_{\pm 0.39}$	77.97 <sub>±</sub>
CIFAR10-S	10 50 100	$23.60 \pm 0.32 \\ 36.46 \pm 0.49 \\ 39.34 \pm 0.56$	$37.88_{\pm 0.27}$ $45.02_{\pm 0.44}$ $48.11_{\pm 0.63}$	$\begin{array}{c} 45.17_{\pm 0.46} \\ 58.84_{\pm 0.23} \\ 61.33_{\pm 0.37} \end{array}$	$\begin{array}{c} 37.88 \pm 0.76 \\ 41.28 \pm 0.80 \\ 42.73 \pm 0.73 \end{array}$	$\begin{array}{c} 41.82_{\pm 0.70} \\ 49.26_{\pm 0.42} \\ 51.74_{\pm 0.52} \end{array}$	$\begin{array}{c} 48.30_{\pm 0.79} \\ 47.26_{\pm 0.73} \\ 47.27_{\pm 1.66} \end{array}$	$56.40_{\pm 0.37}$ $57.84_{\pm 0.24}$ $56.98_{\pm 0.85}$	$\begin{array}{c} 55.09 {\pm} 0.43 \\ 57.59 {\pm} 0.93 \\ 57.14 {\pm} 0.26 \end{array}$	$58.40_{\pm 0.24}$ $61.85_{\pm 0.33}$ $62.70_{\pm 0.37}$	69.78 <sub>±</sub>
CelebA	10 50 100	$54.51 \pm 1.63$ $55.99 \pm 1.23$ $60.62 \pm 0.71$	$61.79_{\pm 0.82}$ $64.61_{\pm 0.25}$ $65.13_{\pm 0.30}$	$64.37_{\pm 0.31}$ $68.50_{\pm 0.58}$ $68.84_{\pm 0.32}$	$57.19_{\pm 2.31}$ $60.16_{\pm 1.57}$ $62.53_{\pm 2.42}$	$57.63_{\pm 1.49}$ $59.89_{\pm 1.11}$ $61.89_{\pm 1.91}$	$61.49_{\pm 0.57}$ $60.75_{\pm 0.37}$ $64.04_{\pm 0.68}$	$63.54_{\pm 0.73}$ $66.89_{\pm 0.28}$ $67.24_{\pm 0.58}$	$\begin{array}{c} 64.38_{\pm 0.33} \\ 64.62_{\pm 0.44} \\ 62.58_{\pm 0.45} \end{array}$	$66.26_{\pm 0.18}$ $68.26_{\pm 0.37}$ $64.12_{\pm 0.28}$	74.09±

Method	Cross		DM			DM+FairDD	)
Wiethou	arch.	DEOM	DEOA	Acc.	DEOM	DEOA	Ace
	ConvNet	$100.0 \pm 0.00$	$91.68 \pm 1.67$	$56.84_{\pm 1.92}$	$10.05 \pm 0.75$	$5.46 \pm 0.40$	$96.58\pm$
C-MNIST	AlexNet	$100.0 \pm 0.00$	$98.82 \pm 1.63$	$44.02 {\pm} 3.34$	$10.35 {\pm} 0.60$	$6.16 \pm 0.39$	$96.12 \pm$
(FG)	VGG11	$99.70_{\pm 0.02}$	$70.73_{\pm 3.44}$	$75.22_{\pm 1.41}$	$9.55 \pm 0.66$	$5.39_{\pm 0.36}$	$96.80_{\pm}$
(10)	ResNet18	$100.0 \pm 0.00$	$96.00_{\pm 1.07}$	$52.05 \pm 1.93$	$8.40 \pm 0.50$	$4.63 \pm 0.27$	97.13 <sub>±</sub>
	Mean	99.93	89.31	57.03	9.59	5.41	96.0
	ConvNet	$ 100.0_{\pm 0.00} $	$99.71_{\pm 0.10}$	$36.27 \pm 0.72$	$24.50 \pm 1.19$	$14.47_{\pm 0.66}$	79.07±
C EMNIST	AlexNet	$100.0_{\pm 0.00}$	$99.75_{\pm 0.12}$	$22.72_{\pm 1.60}$	$20.60 \pm 1.03$	$14.11 \pm 0.85$	$76.14 \pm$
(BG)	VGG11	$100.0 \pm 0.00$	$97.77_{\pm 0.98}$	$43.11 \pm 0.79$	$21.60 \pm 0.76$	$14.36 \pm 0.91$	$78.57$ $\pm$
(DC)	ResNet18	$100.0 \pm 0.00$	$99.78 \pm 0.09$	$23.37_{\pm 1.11}$	$22.50 \pm 1.09$	$14.96 \pm 0.81$	75.21 <sub>±</sub>
	Mean	100.0	99.25	31.37	22.30	14.73	77.2
	ConvNet	$ 75.13_{\pm 1.24} $	$55.70_{\pm 1.11}$	$45.02 \pm 0.44$	$18.28 \pm 0.78$	$7.35 \pm 0.59$	58.84 <sub>±</sub>
	AlexNet	$75.30 \pm 0.90$	$52.57 \pm 1.29$	$36.09 \pm 0.51$	$15.84 \pm 1.11$	$5.12 \pm 0.28$	49.16 <sub>±</sub>
CIFAR10-S	VGG11	$61.48 \pm 1.91$	$44.05 \pm 1.40$	$43.23 {\pm} 0.55$	$11.51 \pm 0.77$	$4.16 \pm 0.32$	52.65±
	ResNet18	$76.23 \pm 0.72$	$54.35 \pm 0.99$	$38.03 \pm 0.51$	$16.44 \pm 0.98$	$5.14 \pm 0.50$	50.93±
	Mean	72.04	51.67	40.59	15.27	5.44	52.9
	ConvNet	$ 40.26 \pm 0.94 $	$38.81 \pm 0.70$	64.61±0.25	$14.08 \pm 0.32$	$9.87_{\pm 1.00}$	68.50 <sub>±</sub>
	ConvNet AlexNet	$ 40.26 \pm 0.94 $ $ 32.51 \pm 0.85 $	$38.81_{\pm 0.70}$ $31.62_{\pm 1.03}$	$64.61_{\pm 0.25}$ $63.10_{\pm 0.55}$	$14.08 \pm 0.32$ $9.38 \pm 0.86$	$9.87_{\pm 1.00}$ $5.75_{\pm 0.76}$	68.50± 64.24±
CelebA	ConvNet AlexNet VGG11	$\begin{array}{ } 40.26_{\pm 0.94} \\ 32.51_{\pm 0.85} \\ 26.03_{\pm 2.10} \end{array}$	$\begin{array}{c} 38.81 {\pm} 0.70 \\ 31.62 {\pm} 1.03 \\ 24.63 {\pm} 1.97 \end{array}$	$\begin{array}{c} 64.61 {\pm} 0.25 \\ 63.10 {\pm} 0.55 \\ 61.57 {\pm} 0.79 \end{array}$	$\begin{array}{c} 14.08 {\pm} 0.32 \\ 9.38 {\pm} 0.86 \\ 8.95 {\pm} 0.97 \end{array}$	$\begin{array}{c} 9.87_{\pm 1.00} \\ 5.75_{\pm 0.76} \\ 6.32_{\pm 1.10} \end{array}$	$68.50_{\pm}$ $64.24_{\pm}$ $62.05_{\pm}$
CelebA	ConvNet AlexNet VGG11 ResNet18	$\begin{array}{c} 40.26 \pm 0.94 \\ 32.51 \pm 0.85 \\ 26.03 \pm 2.10 \\ 25.60 \pm 1.87 \end{array}$	$\begin{array}{c} 38.81 {\scriptstyle \pm 0.70} \\ 31.62 {\scriptstyle \pm 1.03} \\ 24.63 {\scriptstyle \pm 1.97} \\ 24.93 {\scriptstyle \pm 1.75} \end{array}$	$\begin{array}{c} 64.61 {\pm} 0.25 \\ 63.10 {\pm} 0.55 \\ 61.57 {\pm} 0.79 \\ 60.32 {\pm} 0.82 \end{array}$	$\begin{array}{c} 14.08 {\pm} 0.32 \\ 9.38 {\pm} 0.86 \\ 8.95 {\pm} 0.97 \\ 6.72 {\pm} 0.81 \end{array}$	$\begin{array}{c} 9.87_{\pm 1.00} \\ 5.75_{\pm 0.76} \\ 6.32_{\pm 1.10} \\ 4.29_{\pm 0.67} \end{array}$	$\begin{array}{c} 68.50_{\pm} \\ 64.24_{\pm} \\ 62.05_{\pm} \\ 61.80_{\pm} \end{array}$

**Table 23:** Cross-arch. comparison. The results are reported in the format of mean  $\pm$  standard deviation.

Table 24: Ablation on BR at IPC = 50. The results are reported in the format of mean  $\pm$  standard deviation.

Methods			DM			DM+FairDD	)
Dataset	BR	DEOM	$\overline{\text{DEO}_{A}}$	Acc.	$\overline{\text{DEO}_M}$	$\overline{\text{DEO}_A}$	Acc.
C-MNIST (FG)	0.85 0.90 0.95	$\begin{array}{c} 99.54_{\pm 0.37} \\ 100.0_{\pm 0.00} \\ 100.0_{\pm 0.00} \end{array}$	$70.13_{\pm 1.69} \\91.68_{\pm 1.67} \\100.0_{\pm 0.00}$	$\begin{array}{c} 76.24_{\pm 1.68} \\ 56.84_{\pm 1.92} \\ 33.73_{\pm 1.08} \end{array}$	$\begin{array}{c} 10.13 {\scriptstyle \pm 0.75} \\ 10.05 {\scriptstyle \pm 0.75} \\ 10.30 {\scriptstyle \pm 0.74} \end{array}$	$5.20_{\pm 0.33} \\ 5.46_{\pm 0.40} \\ 5.84_{\pm 0.39}$	$96.62_{\pm 0.1} \\ 96.58_{\pm 0.0} \\ 96.05_{\pm 0.1}$
C-FMNIST (BG)	0.85 0.90 0.95	$100.0_{\pm 0.00} \\ 100.0_{\pm 0.00} \\ 100.0_{\pm 0.00}$	$\begin{array}{c} 95.54_{\pm 0.79} \\ 99.71_{\pm 0.10} \\ 99.79_{\pm 0.08} \end{array}$	$\begin{array}{c} 46.14_{\pm 0.93} \\ 36.27_{\pm 0.72} \\ 26.30_{\pm 0.44} \end{array}$	$\begin{array}{c} 23.75_{\pm 1.58} \\ 24.50_{\pm 1.19} \\ 29.15_{\pm 1.30} \end{array}$	$\begin{array}{c} 13.85_{\pm 0.82} \\ 14.47_{\pm 0.66} \\ 17.72_{\pm 0.90} \end{array}$	$79.61_{\pm 0.1} \\ 79.07_{\pm 0.2} \\ 78.46_{\pm 0.2}$
CIFAR10-S	0.85 0.90 0.95	$71.75_{\pm 0.90} \\ 75.13_{\pm 1.24} \\ 75.43_{\pm 1.28}$	$\begin{array}{c} 50.11 {\scriptstyle \pm 0.70} \\ 55.70 {\scriptstyle \pm 1.11} \\ 58.58 {\scriptstyle \pm 0.82} \end{array}$	$\begin{array}{c} 46.99 _{\pm 0.43} \\ 45.02 _{\pm 0.44} \\ 43.56 _{\pm 0.38} \end{array}$	$\begin{array}{c} 16.44_{\pm 0.88} \\ 18.28_{\pm 0.78} \\ 17.49_{\pm 1.26} \end{array}$	$\begin{array}{c} 6.58_{\pm 0.76} \\ 7.35_{\pm 0.59} \\ 7.10_{\pm 0.91} \end{array}$	$59.12_{\pm 0.32}$ $58.84_{\pm 0.22}$ $58.18_{\pm 0.34}$

Table 25: Ablation on initialization at IPC = 50. The results are reported in the format of mean  $\pm$ standard deviation.

Methods	Init		DM			DM+FairDD	)
Dataset	11111.	DEOM	DEOA	Acc.	DEOM	DEOA	Ac
C-MNIST	Random	$ 100.0_{\pm 0.00} $	$91.68 \pm 1.67$ 99.64 + 0.80	$56.84_{\pm 1.92}$	$10.05 \pm 0.75$	$5.46_{\pm 0.40}$	96.58±
(FG)	Hybrid	$100.0\pm0.00$ $100.0\pm0.00$	99.06 $\pm 0.74$	$39.97 \pm 1.77$	$9.03 \pm 0.63$ $9.03 \pm 0.62$	$5.23 \pm 0.29$ $5.33 \pm 0.25$	$96.27_{\pm}$
C-FMNIST	Random	$ 100.0_{\pm 0.00} $	$99.71 \pm 0.10$	$36.27 \pm 0.72$	$24.50 \pm 1.19$	$14.47 \pm 0.66$	79.07±
(BG)	Noise Hybrid	$100.0 \pm 0.00$ $100.0 \pm 0.00$	$99.67 \pm 0.07$ $99.68 \pm 0.07$	$22.92 \pm 0.78$ $26.38 \pm 0.80$	$23.00 \pm 1.93$ $21.45 \pm 1.32$	$14.40 \pm 0.63$ $14.34 \pm 0.89$	/8.84 <sub>∃</sub> 79.19 <sub>∃</sub>
	Random	$ 75.13_{\pm 1.24} $	55.70 <sub>±1.11</sub>	$45.02_{\pm 0.44}$	$18.28 \pm 0.78$	$7.35_{\pm 0.59}$	58.84 <sub>±</sub>
CIFAR10-S	Noise Hybrid	$55.28 \pm 1.57$ $65.59 \pm 1.38$	$37.26 \pm 0.63$ $46.18 \pm 0.84$	$46.97 \pm 0.29$ $45.16 \pm 0.40$	$16.15 \pm 0.62$ $17.30 \pm 1.26$	$6.13_{\pm 0.50}$ $6.71_{\pm 0.33}$	56.41 <sub>±</sub> 56.78 <sub>±</sub>