

---

# [Re] Deep Fair Clustering for Visual Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 The authors propose a novel method for Deep Fair Clustering (DFC), combining existing frameworks for fair clustering—which typically have difficulty with high-dimensional large-scale data—with previous work on deep clustering—which typically has difficulty with fairness. Our reproducibility work targets the central claim that DFC learns fair representations with minimal utility loss and obtains superior results on both fairness and accuracy.

### 7 **Methodology**

8 We used the code repository made available by the authors and extended it to include support for pretraining, different datasets and comparative methods. We compare the DFC method against Deep Embedded Clustering (DEC) (11), which implements a comparable deep clustering method without fairness constraints, on the same four datasets (obtained from MNIST (7), USPS (6), MTFL (13) and Office-31 (9)) and fairness metrics that were used in the paper. We select one dataset (MNIST-USPS) for additional experiments aimed at validating the contribution of individual components of the DFC towards fairness. All experiments were run on a GeForce 1080Ti GPU. Hyperparameter optimization was performed using the Weights & Biases Sweeps feature (1).

### 15 **Results**

16 On the selected dataset, we reproduced accuracy to within 2% of reported value, normalized mutual information (NMI) and entropy to within 1%, and balance to within 5%. Our DFC method outperformed our DEC method on all accuracy and fairness metrics. We reproduced the accuracy of the non-digit datasets to within 1% (Office-31) and 7% (MTFL) but failed to obtain similar results for balance.

### 20 **What was easy**

21 We found no major challenges reproducing the provided code in as far as we used the author’s provided pretrained models and the selected dataset (MNIST-USPS) used in the code.

### 23 **What was difficult**

24 Extending the code for the non-digit datasets was a challenge, as some hyperparameter settings and architecture details were difficult to infer from the paper. We found that performance was sensitive to small changes in the implementation and training of the encoders. Consequently, we ran into time and resource constraints when trying to reproduce all results for these datasets, due to the large number of models that required pretraining.

### 28 **Communication with original authors**

29 We had helpful one-off contact with the authors to verify hyperparameter settings.

## 30 1 Introduction

31 Machine learning (ML) is increasingly used in high-stake decision-making where the data contains sensitive attributes,  
32 such as gender, race or socioeconomic background. Examples include college admission, loan approval and bail/parole  
33 judgements. Such ML algorithms are vulnerable to bias and unfairness (4), and extensive literature has brought attention  
34 to the various challenges and inherent trade-offs in this field (5)(14).

35 Clustering is an important aspect of many such ML applications. To illustrate what it might mean for a clustering  
36 assignment to be unfair, we can consider its use in feature engineering, e.g., to label data in an unsupervised setting to  
37 increase its expressive power (2). In some cases, if the inherent structure of the training data differs between subgroups,  
38 standard clustering likely leads to partition assignments that correlate significantly with sensitive attributes. Often, these  
39 are controversial correlations that do not reflect true causal mechanisms. This could facilitate subsequent discrimination  
40 indirectly based on sensitive attributes, either knowingly or unknowingly.

41 In short, fair clustering is an ongoing, important, and complex field that still faces many challenges. The paper 'Deep  
42 Fair Clustering for Visual Learning' addresses various such challenges. Moreover, their method for fair clustering  
43 outperforms competitive fair clustering methods as well as competitive deep clustering methods. This shows that, if  
44 theoretically sound, their method is able to impose fairness constraints without significantly compromising clustering  
45 accuracy.

## 46 2 Scope of reproducibility

47 The paper aims to tackle the problem of fair clustering of high-dimensional data by introducing a Deep Fair Clustering  
48 (DFC) method suitable for image data. Their notion of fairness extends the demographic equipartition criterion  
49 and requires that the clustering assignment is independent of the protected subgroup membership. To achieve this  
50 independence, they leverage the feature representation encoding from deep methods and train a model to filter out  
51 sensitive attributes from the representations. Simultaneously, DFC enforces fairness constraints and optimizes clustering  
52 performance, notably with a minimal trade-off between the two. The central claims are summed up as follows:

- 53 • The first claim states that DFC achieves fair clustering of large-scale and high-dimensional visual data by  
54 learning fair representations of the input. To support this claim, we train DFC and DEC on the MNIST-USPS  
55 dataset and compare the learned representations in figure 3. The first row of table 1 shows the corresponding  
56 accuracy and fairness metrics of the final performance of both methods.
- 57 • The second claim regards the validity and effectiveness of the proposed minimax optimization formulation. In  
58 particular, we support this claim by recognizing DEC as DFC without the minimax optimization formulation  
59 (in particular, without the fairness components 4, 5 and the separate clustering of subgroups, see figure 1).  
60 To our best knowledge, this corresponds to the implementations in (11) and (8). The claim is supported by  
61 comparing DFC performance with DEC performance on accuracy and fairness metrics in table 1. In addition,  
62 figure 4 shows the performance of the discriminator during DFC, which provides some insight in the extend to  
63 which DFC achieves 'masked' representations during training.
- 64 • The third claim states that DFC shows superior performance on four real-world visual datasets (see 3.1). We  
65 aim to show this by running experiments on the same four datasets with both DEC and DFC and compare  
66 performance in table 1.

## 67 3 Methodology

68 Each claim requires us to implement the DFC method and the DEC method and train it on the MNIST-USPS dataset.  
69 From here, we address each claim separately in three sets of experiments. For the first, we visualize the representations  
70 using dimensionality reduction methods and compare the obtained clusters. For the second, we compare the DEC and  
71 DFC on accuracy, NMI, Balance and Entropy metrics. Moreover, we monitor discriminator accuracy for DFC. For the  
72 third, we run our experiments with the three other datasets that are used in the paper and compare results on the above  
73 mentioned four metrics.

74 Before discussing the details, we briefly explain the DFC method as we implemented it based on the authors' description.  
75 The method has several components that are trained simultaneously in a minimax optimization scheme. An overview of  
76 the method is shown in figure 1. The task is to cluster datapoints  $X$  with sensitive categorical attribute  $G$  belonging to  
77 protected subgroup  $G(X) \in [M]$  into  $K$  clusters. A feature encoder  $\mathcal{F}(X)$  transforms the data  $X$  into representations  
78  $Z$ . As the figure shows, each subgroup is clustered separately and requires a separate, pretrained encoder. The clustering  
79 assignment  $P = \mathcal{A}(Z)$ , taken from the Deep Embedded Clustering method (DEC), creates a soft assignment  $P$

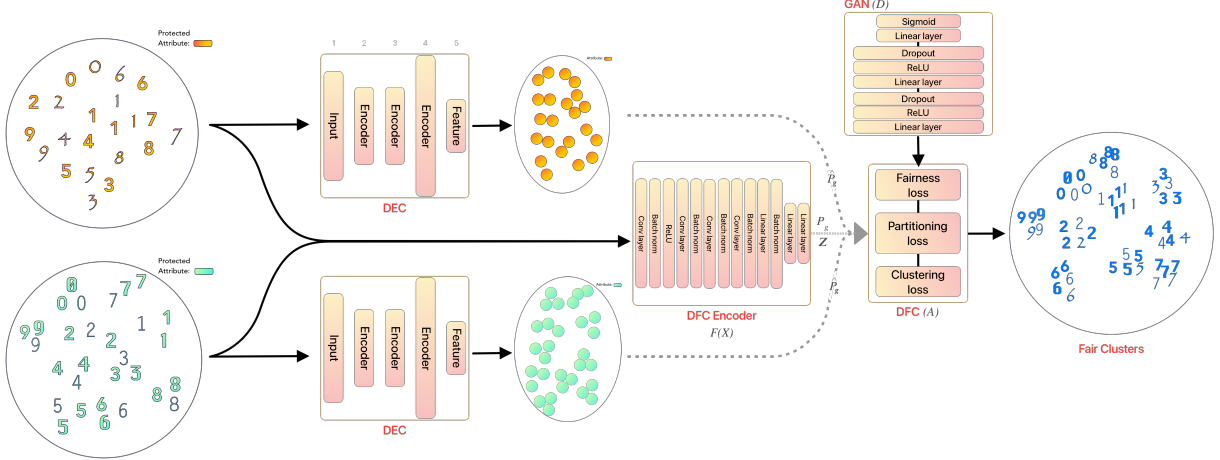


Figure 1: Overview of Deep Fair Clustering. The orange and green colors represent the protected subgroups.

80 reflecting the probability of assigning datapoints to each cluster, once for each protected subgroup, and once for the full  
81 dataset.

82 A discriminator  $\mathcal{D}$  aims to reconstruct the protected subgroup membership based on the soft assignment  $P$

83 **Encoder** Central to the DFC approach is the encoding of fair representations. Transforming the input data using  
84 feature encoding is important, in the first place, because (deep) clustering at pixel level does not often yield good results.  
85 In DFC, the encoded representations are at the same time an important step towards fairness. The authors note that  
86 "fair representations are achieved when the discriminator cannot distinguish between representations from different  
87 protected subgroups." To this end, the encoder is finetuned simultaneously with the rest of the clustering algorithm.

88 **Discriminator** A discriminator is constructed to monitor and guide the fairness of representations  $Z$ .

89 **Deep Embedded Clustering** For the clustering assignment, the authors closely follow the DEC method. As in (11),  
90 the cluster assignment is learned using Stochastic Gradient Descent (SGD), repeating the following two processes:

- 91 1. A soft assignment of the datapoints to cluster centroids, using the Student's t-distribution as a kernel to measure  
92 the similarity between embedded point  $z$  and centroid  $c_k$ :

$$p_k = \frac{(1 + \|z - c_k\|^2/\alpha)^{-\frac{\alpha+1}{2}}}{\sum_{k' \in [K]} (1 + \|z - c_{k'}\|^2/\alpha)^{-\frac{\alpha+1}{2}}}, \quad (1)$$

93 where  $p_k$  is the probability of a datapoint belonging to cluster center  $c_k$  and  $\alpha$  the degrees of freedom of the  
94 Student's t-distribution.

- 95 2. The obtained distribution  $P$  of soft assignments is matched to an auxiliary distribution  $Q$ :

$$q_k = \frac{(p_k)^2 / \sum_{x \in X_g} p_k}{\sum_{k' \in [K]} ((p_{k'})^2 / \sum_{x \in X_g} p_{k'})}, \quad (2)$$

96 which is calculated for each protected subgroup separately. Finally, the clustering is optimized by minimizing  
97 the KL-divergence loss:

$$\mathcal{L}_c := KL(P||Q) = \sum_{g \in [M]} \sum_{x \in X_g} \sum_{k \in [K]} p_k \log \frac{p_k}{q_k}. \quad (3)$$

98 This ensures that the soft assignment matches the target distribution as closely as possible.

99 **Deep Fair Clustering** DFC implements DEC as above, but adds two loss functions to impose the fairness constraint  
100 while promoting cluster utility. To encourage fair partitions according to the demographic parity, the authors introduce  
101 a *fairness adversarial loss*:

$$\mathcal{L}_f := \mathcal{L}(\mathcal{D} \circ \mathcal{A} \circ \mathcal{F}(X), G), \quad (4)$$

102 Where  $\mathcal{L}$  is the cross-entropy loss function and  $\circ$  denotes function composition, in particular the result of the encoder  
103  $\mathcal{F}(X)$ , the clustering assignment  $\mathcal{A}$  and the discriminator  $\mathcal{D}$  applied to the data in sequence. The discriminator  $\mathcal{D}$  is

104 implemented as a Generative Adversarial Network (GAN) (see figure 1) with the exact inverse architecture of the used  
 105 encoder. It monitors how well the DFC is able to reconstruct the protected subgroup based on the predicted cluster  
 106 probabilities of a datapoint. With respect to the *fairness adversarial loss*,  $\mathcal{D}$  is maximized with the probability of  
 107 assigning the correct membership label for each sample. Simultaneously, clustering  $\mathcal{A}$  and representations  $\mathcal{F}(X)$  are  
 108 trained to maximally confuse  $\mathcal{D}$ .

109 To increase the clustering performance under the fairness constraint, DFC adds a *structure preservation loss*:

$$\mathcal{L}_s = \sum_{g \in [M]} \left\| \hat{P}_g \hat{P}_g^\top - P_g P_g^\top \right\|^2, \quad (5)$$

110 where  $\hat{P}_g$  is the soft assignment for a datapoint in protected subgroup  $g$ , and  $P_g$  is the result for subgroup  $g$  when the  
 111 clustering was performed over the complete dataset. In other words, based on theoretical considerations as described in  
 112 the original paper (8), DFC expects local and global subgroup clustering partition to be similar.

113 The final learning objective is as follows:

$$\max_{\mathcal{F}, \mathcal{A}} \alpha_f \mathcal{L}_f - \alpha_s \mathcal{L}_s - \mathcal{L}_c \quad (6)$$

$$\min_D \alpha_f \mathcal{L}_f \quad (7)$$

114 Note how the fairness adversarial loss (4) is simultaneously maximized w.r.t.  $\mathcal{F}$  and  $\mathcal{A}$  (6) to minimize performance of  
 115 the discriminator as a result of clustering, and minimized w.r.t.  $\mathcal{D}$  (7) to optimize its parameters and obtain a saddle-point  
 116 solution. The values for  $\alpha$  are 1 by default. Following the paper, for the office-31 dataset  $\alpha_s$  is recalculated as  
 117  $(512/128)^2 \cdot (31/10)$  to adjust for difference in batch size and cluster number.

### 118 3.1 Datasets

119 We perform our experiments on the same four datasets as in the original paper. To test claim 1 and 2, we restrict  
 120 ourselves to the *MNIST-USPS* dataset, obtained from combining the *MNIST* dataset (7) and the *USPS* dataset (6).  
 121 From both, only the train images are taken, 60,000 and 7,261 respectively. The *Color Reverse MNIST* is obtained by  
 122 reversing the black and white pixels in the MNIST dataset and combining the normal and the reversed version. The  
 123 *The Multi-task Facial Landmark (MTFL)* dataset (13) consist of facial recognition pictures, from which 1000 images  
 124 with and 1000 without glasses are sampled randomly. Finally, the *The Office-31* dataset (9) consists of images from 31  
 different categories collected from two distinct domains: Amazon and Webcam. For important statistics, see figure 2.

Dataset	Type	# Instances	# Classes	# Dim.	Sensitive Attribute	Max of Bal.	Max of Ent.
MNIST-USPS	digital	67,291	10	1,024	source of digital	0.120	2.303
Color Reserve MNIST	digital	120,000	10	1,024	original or reversed	1.000	2.303
MTFL	face	2,000	2	50,176	with or w/o glasses	1.000	0.693
Office-31(A&W)	object	3,612	31	50,176	domain source	0.273	3.434

Figure 2: Characteristics of each of the four datasets, taken from the original paper (8).

125

#### 126 3.1.1 Hyperparameters

127 The DFC model is dependent on 5 pretrained models; a VAE and a DEC for each subgroup and a VAE for the whole  
 128 dataset. Since the optimal hyperparameters of the DFC are likely to depend on the hyperparameter choices of the  
 129 pretrained models, a complete search would be expensive. To make hyperparameter search manageable we make the  
 130 assumption that the optimal hyperparameters are independent of the input data and the pretrained models it uses. This  
 131 assumption enables to search the hyperparameter space of the DFC, DEC and VAE independently. Note that the VAE  
 132 only needs to be pretrained and optimized for the MNIST-USPS and Reverse MNIST datasets, the other datasets use a  
 133 pretrained ResNET50 encoder.

134 The hyperparameter search for the DEC and the VAE focused on finding the best combination of the learning rate,  
 135 batch size and the number of epochs. These searches were performed with Bayesian optimization and contained  
 136 respectively 32 and 18 experiments on MNIST. The VAE Bayesian search was set to maximize the accuracy of the  
 137 KMeans clustering afterwards, which is also used in the training of the DFC as initial centroids. The goal of the DEC  
 138 was to optimize its clustering accuracy.

139 The DEC and VAE were pretrained using the optimal parameters on MNIST and USPS to facilitate the hyperparameter  
140 search of the DFC. The focus was on finding the optimal learning rate, batch size, and number of epochs. A random  
141 grid search was used because it is not clear which metric should be optimized. The best performing set of parameters  
142 was chosen based on the Accuracy, NMI, entropy and balance.

143 The range of the hyperparameter space was inspired by hyperparameters mentioned by the authors in the paper, email  
144 contact and a DEC MNIST example<sup>1</sup>.

145 The hyperparameters of the KMeans clustering were set to the same values the authors used. TSNE hyperparameters  
146 were set to defaults as can be seen in Appendix C, which contains a complete list of all hyperparameters and overviews  
147 of the hyperparameter searches.

## 148 **4 Experimental setup and implementation**

149 All experiments were run on cuda enabled computers with seeds set to 2019 before training each model, unless otherwise  
150 specified. Our code containing the experiments is available at our GitHub repo<sup>2</sup>. Our code is strongly based on the  
151 implementation provided by the authors<sup>3</sup>. We have found that their code is a clear reflection of the paper. However, it  
152 was necessary to extend their implementation to provide support for pretraining, other datasets and other methods.

153 The DFC architectures described in the paper relies on several models joint together. Roughly, it includes an encoder, a  
154 clustering assignment and a discriminator. Some experiments require slightly different model combinations, which is  
155 always made explicit.

### 156 **4.1 Pretrained Encoder**

157 The type of encoder can be chosen independently from the DFC. The authors use different architectures based on the  
158 complexity of the dataset; for the digit based datasets a standard Variational Autoencoder is used, for the other datasets  
159 a pretrained ResNet50.

160 The standard VAE is based on four blocks of a 2D convolutional layer, batch normalization and a ReLU activation layer.  
161 This is scaled down in 2 steps with fully connected layers to the 64 dimensional hidden space. Since the decoder was  
162 missing from the author’s implementation, some guesses had to be made about the exact implementation, we tried to  
163 follow conventions as much as possible. A standard KL-Divergence and MSE loss were used for pretraining. The exact  
164 architecture can be found in Appendix A.

165 For MTFL and Office-31, considering their complex and high resolution images, the authors use ImageNet-pretrained  
166 ResNet50 (3). The authors did not mention how the ResNet is finetuned exactly. Due to the relatively small size of the  
167 datasets used for the experiments, it was decided to only finetune the last layer during training the DFC.

### 168 **4.2 Pretrained DEC**

169 The structural preservation loss, discussed before 5, is calculated using two separate DEC models. These models could  
170 be seen as the golden standard for each subgroup. Since they are equally weighted, this ensures that the protected  
171 subgroup is underrepresented in the loss. These DEC’s also require an encoder to transform the images. The same  
172 encoder training methods and architectures are used for all encoders. The DEC’s were implemented as a stripped down  
173 DFC with three modifications. The fairness and the structural preservation loss were removed. Lastly, the clustering loss  
174 is calculated over the entire batch at once, instead of individually per subgroup. This results in an exact implementation  
175 of the original DEC paper (11).

### 176 **4.3 Centroid Initialization**

177 The deep clustering methods DEC and DFC both benefit from properly initializing the Student’s t-distribution cluster  
178 definitions. This is done with a simple KMeans clustering approach. KMeans is trained on the entire dataset that is used  
179 by the clustering method afterwards. The final centroids learned are copied to be the initial cluster definitions.

---

<sup>1</sup><https://github.com/vlukiyarov/pt-dec/blob/master/examples/mnist/mnist.py>

<sup>2</sup><https://github.com/Joppewouts/Reproducing-Deep-Fair-Clustering>

<sup>3</sup><https://github.com/brandeis-machine-learning/DeepFairClustering>

## 180 4.4 DFC

181 The DFC consists of a pretrained encoder, a discriminator module and cluster assignment module. The weights of these  
182 modules are updated with the Adam optimizer, using an inverse exponential learning rate scheduler. The discriminator  
183 has a standard architecture, details are attached in Appendix B. As described by the authors, the learning rate of the  
184 discriminator uses a multiplication factor of 10.

### 185 4.4.1 Clustering Assignment

186 The fairness constraints and learning objectives of DFC rely on clustering being performed on the subgroups separately.  
187 Since these subgroups are homogeneous with respect to sensitive attributes, they can be clustered irrespective of  
188 fairness-considerations. Therefore, the same model can be used as in the Deep Embedding clustering. Accordingly, the  
189 authors use the DEC clustering for the subgroups separately, as described in section 3. For full details see also (11).

## 190 4.5 Competitive methods

191 Due to time constraints, we chose to only implement one competitive method. DEC was a compelling candidate because  
192 it allows to compare fair deep clustering with ordinary deep clustering.

## 193 4.6 Metrics

194 Four metrics are used to evaluate performance; accuracy, normalized mutual information, balance and entropy. These  
195 are the same metrics as in the original paper, which includes detailed explanations. Furthermore, for evaluating the  
196 fairness two more metrics are implemented. The accuracy of the discriminator is computed to function as a proxy for the  
197 fairness of the cluster assignments. The probabilities outputted by the discriminator are converted to subgroup classes by  
198 thresholding at 0.5. An accuracy of 50% indicates fully masked representations, and therefore fair cluster assignments.  
199 Secondly, TSNE (10) is implemented to visualize the latent space learned by the encoder. In an ideal fair representation,  
200 the representation of the two subgroups should be indistinguishable. For example, the digit 5 of MNIST and the digit 5  
201 of USPS both encoded as the character '5' would be ideal, since the representations are indistinguishable. The goal of a  
202 TSNE mapping is that distant samples in the high dimensional space are also far away from each other in the lower  
203 dimension, and vice versa. Therefore, in a visualization of a fair latent space, it should be hard to distinguish clusters of  
204 protected subgroups. All results are calculated over 5 different seeds: 0, 1, 42, 2019 and 2020. However, it was not  
205 possible to provide the variance for the results that require a ResNet encoder due to time and resource limitations.

### 206 4.6.1 Computational requirements

207 We trained all models on an 8.25MB GeForce 1080Ti GPU. Pretraining VAE encoders for each dataset took  $(5 \times) \pm$   
208 45min, pretraining the DEC models  $(7 \times) \pm$  1hr. The sweep for optimal parameter search took  $(20 \times)$  1hr and running  
209 the DFC and DEC models for all four datasets on 5 random seeds took  $(8 \times) \pm$  1hr GPU time. The total experiments  
210 required a total of  $\pm 40$ hr GPU time.

## 211 5 Results

212 Our results support the first claim that DFC learns fair representations that do not cause much loss in utility and the  
213 second claim that the minimax optimization formulation contributes to clustering that is both fair and accurate. We  
214 reproduced accuracy on the MNIST-USPS dataset to within 2% of reported value, normalized mutual information  
215 (NMI) and Entropy to within 1%, and balance to within 5%. Our DFC method outperformed our DEC method on  
216 all accuracy and fairness metrics. We reproduced the accuracy of the Office-31 datasets to within 1% and the MTFI  
217 dataset to within 7%, but failed to obtain similarly good results for the fairness metrics.

### 218 5.1 Claim 1: Learning fair representations for the MNIST-USPS dataset

219 Figure 3 shows the learned representation of the DFC method (top) and the DEC method (bottom). Protected subgroup  
220 membership is indicated by color for MNIST (red) and USPS (blue). Monitoring the DFC discriminator performance  
221 shows that the accuracy for recovering the protected subgroup drops from 60% at the first 1000 runs to 50% after 2000  
222 iterations. The corresponding accuracy and fairness metrics are shown in table 1. We obtain an accuracy within 2% of  
223 the accuracy reported in the paper. NMI and Entropy scores fell within 1% of the original findings, and balance was  
224 5.5% lower in our experiments.

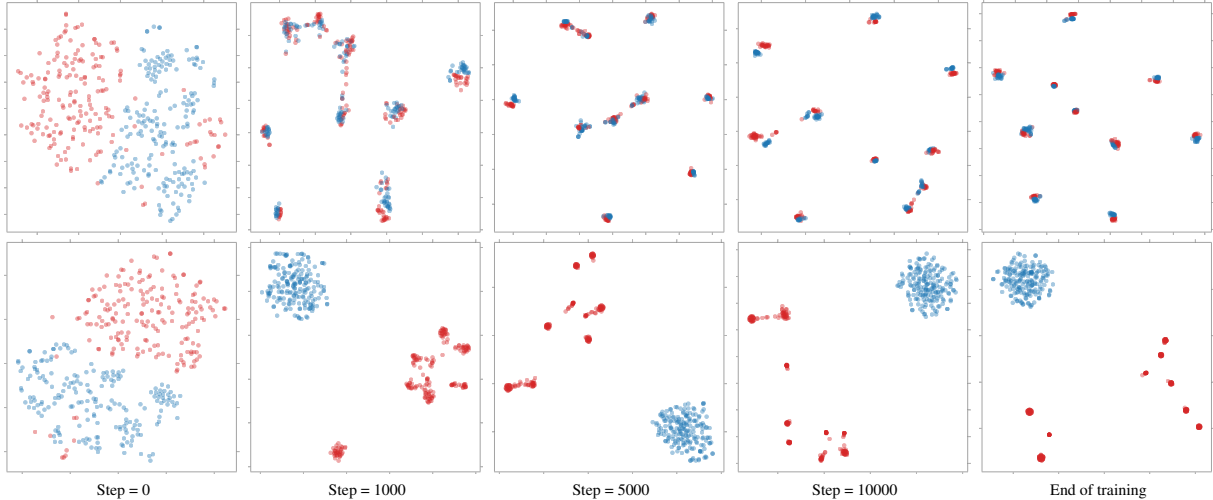


Figure 3: TSNE visualizations of the latent space learned by the encoder on MNIST-USPS with the final models. The rows respectively show the representations for the DFC and for the DEC. The DFC learns representations of the images which results nicely in 10 clusters in the visualized space. The DEC learns 9 distinct representation for MNIST digits and does not seem to separate USPS digits at all.

## 225 5.2 Claim 2: Validity of minimax optimization formalization

226 The first two rows of table 1 show performance on all  
 227 four metrics for the DFC and DEC methods. Results  
 228 indicate that DFC outperforms DEC both in accuracy  
 229 and in fairness methods. Moreover, the table shows that  
 230 our implementation of the DEC method achieves a 17%  
 231 higher accuracy and performs better than the original  
 232 implementation of the DEC on all other metrics except  
 233 balance.

234 To further investigate how subcomponents of DFC con-  
 235 tribute towards performance and fairness, we plot the  
 236 decoder accuracy in figure 4. Interestingly, the plot shows  
 237 that DFC initially achieves perfectly masked representa-  
 238 tion up to around 8.5k training steps, where  $\mathcal{F}(X)$  and  
 239  $\mathcal{A}$  no longer succeed in fully confusing  $\mathcal{D}$ , presumably  
 240 in favour of cluster favourable representations. This in-  
 241 dicates that some trade-off is still present.

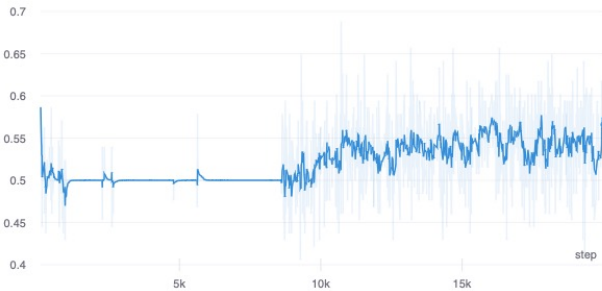


Figure 4: DFC Discriminator accuracy over training steps. An accuracy of 0.5 indicates that representations cannot be discriminated for subgroup membership.

## 242 5.3 Claim 3: Different datasets

243 The results for all datasets are can be found in table 1. Our implementation of DFC outperformed DEC on the three  
 244 other datasets and showed similar accuracy as reported in the original paper. Results for entropy were also consistently  
 245 higher than the DEC implementations and similar to the entropy reported in the original paper. The balance, however,  
 246 was much lower than in the original paper, and did not differ significantly from the DEC results.

## 247 6 Discussion

248 For the MNIST-USPS digit datasets, the results indicate that the method is able to learn fair representations and clusters.  
 249 The entropy results indicate equal fair partitions, the TSNE visualizations show good clusters and the discriminator is  
 250 not able to distinguish between subgroups in the cluster assignments.

251 We were partly unable to reproduce its fairness performance on the three other datasets but can think of various  
 252 limitations of our approach that might provide an explanation for this. First, we put notable effort into adjusting  
 253 and training pretrained models for the MNIST-USPS experiments, and performed a hyperparameter search to find

Dataset		Method	Acc	NMI	Balance	Entropy
MNIST-USPS	ours	DFC	0.81 (0.012)	0.80 (0.015)	0.012 (0.001)	2.292 / 2.299 (0.004 / 0.001)
		DEC	0.76 (0.004)	0.73 (0.003)	0.000 (0.000)	2.166 / 2.278 (0.007 / 0.002)
	orig.	DFC	0.83	0.79	0.067	2.301 / 2.265
		DEC	0.59	0.69	0.000	2.082 / 1.735
Reverse MNIST	ours	DFC	0.39 (0.027)	0.37 (0.038)	0.007 (0.001)	2.283 / 2.293 (0.009 / 0.002)
		DEC	0.32 (0.010)	0.32 (0.018)	0.000 (0.000)	1.395 / 1.939 (0.118 / 0.055)
	orig.	DFC	0.58	0.68	0.763	2.294 / 2.301
		DEC	0.40	0.48	0.000	1.774 / 1.384
MTFL	ours	DFC	0.79	0.27	0.001	0.655 / 0.693
		DEC	0.77	0.23	0.001	0.305 / 0.693
	orig.	DFC	0.72	0.19	0.986	0.693 / 0.693
		DEC	0.52	0.03	0.711	0.660 / 0.576
Office 31	ours	DFC	0.68	0.72	0.000	2.689 / 3.393
		DEC	0.59	0.65	0.000	2.300 / 3.079
	orig.	DFC	0.69	0.72	0.117	3.422 / 3.403
		DEC	0.55	0.60	0.000	3.063 / 2.937

Table 1: Cluster performance and fairness results compared to the original paper with the variance over 5 seeds when applicable.

254 optimal settings. Due to time and resource constraints, we were not able to repeat this process for all datasets, which  
255 may have led to different results. Secondly, our use of the ResNet50 encoder differed slightly from that in the paper.  
256 Correspondence with the authors had confirmed that their implementation finetuned all parameters of ResNet50.  
257 However, for us GPU memory constraints made it impossible to store all gradients that would be required. Therefore,  
258 we resolved to only finetune the final layer, based on transfer learning literature (12) and the relatively small sizes of the  
259 datasets. This could partially explain the difference in results.

260 However, evaluating our full results, an additional question about the added fairness constraints and the experiment  
261 setup of this method arose. In the case of MNIST-USPS, optimal clustering solutions for the datasets closely match  
262 optimal fairness according to the constraints. Since the digit datasets are quite nicely balanced in the sense that they  
263 have similar partitions across subgroups, the equal partition definition of fairness is in line with the optimal solution,  
264 and we suspect that the fairness constraints here do not pose as much of a challenge as they might in more imbalanced  
265 datasets. For this reason, we were particularly interested in the results for the non-digit datasets. Further research might  
266 be needed to investigate the theoretical framework in the context of datasets with different inherent fairness challenges.

## 267 6.1 What was easy and/or difficult

268 Based on the implementation made available by the authors it was straightforward to evaluate the method on MNIST-  
269 USPS, using the provided pretrained encoders, DECes and initial centroids. These results matched the results achieved  
270 in the original paper. However, essential code and specifications were missing to extend these results to other datasets  
271 and methods. Examples include the hyperparameters for pretraining all the models, hyperparameters for running on  
272 other datasets, missing specifications for the Decoder of the VAE and the finetuning of ResNet50. This made it hard to  
273 replicate the method for other datasets.

274 Despite the missing parts, the code was helpful for understanding the method in detail. Most importantly, it cleared up  
275 the workings of the structural preservation loss which was not immediately clear from the paper.

276 The large number of different pretrained models needed was another difficulty. The DFC for MNIST-USPS depends on  
277 8 other models: 3 encoders, 2 DECes and 3 KMeans cluster initializations. This, in combination with the fact that the  
278 DFC is sensitive to the quality of the pretrained models, made it hard to debug performance issues that could be caused  
279 by the model itself or any of the previous methods. In addition, it required excellent experiment management to keep  
280 track of which models are trained with which pretrained models and parameters.

## 281 6.2 Communication with original authors

282 We consulted the authors about hyperparameters for pretraining the encoders and which layers of ResNet50 to finetune.  
283 We received a helpful response in which our questions were adequately answered.



284 **References**

- 285 [1] BIEWALD, L. Experiment tracking with weights and biases, 2020. Software available from wandb.ai.
- 286 [2] CHIERICHETTI, F., KUMAR, R., LATTANZI, S., AND VASSILVITSKII, S. Fair clustering through fairlets. In  
287 *Advances in Neural Information Processing Systems* (2017), I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach,  
288 R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30, Curran Associates, Inc.
- 289 [3] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition. In *Proceedings of the*  
290 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.
- 291 [4] JULIA ANGWIN, JEFF LARSON, S. M., AND KIRCHNER, L. Machine bias: There is software used across the  
292 country to predict future criminals. and it’s biased against blacks., 2016.
- 293 [5] KLEINBERG, J., MULLAINATHAN, S., AND RAGHAVAN, M. Inherent trade-offs in the fair determination of risk  
294 scores, 2016.
- 295 [6] LE CUN, Y., MATAN, O., BOSER, B., DENKER, J. S., HENDERSON, D., HOWARD, R. E., HUBBARD,  
296 W., JACKET, L., AND BAIRD, H. S. Handwritten zip code recognition with multilayer networks. In *[1990]*  
297 *Proceedings. 10th International Conference on Pattern Recognition* (1990), vol. 2, IEEE, pp. 35–40.
- 298 [7] LECUN, Y. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/> (1998).
- 299 [8] LI, P., ZHAO, H., AND LIU, H. Deep fair clustering for visual learning. In *Proceedings of the IEEE/CVF*  
300 *Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2020).
- 301 [9] SAENKO, K., KULIS, B., FRITZ, M., AND DARRELL, T. Adapting visual category models to new domains. In  
302 *European conference on computer vision* (2010), Springer, pp. 213–226.
- 303 [10] VAN DER MAATEN, L., AND HINTON, G. Visualizing data using t-sne. *Journal of Machine Learning Research* 9,  
304 86 (2008), 2579–2605.
- 305 [11] XIE, J., GIRSHICK, R., AND FARHADI, A. Unsupervised deep embedding for clustering analysis. In *Proceedings*  
306 *of the 33rd International Conference on International Conference on Machine Learning - Volume 48* (2016),  
307 ICML’16, JMLR.org, p. 478–487.
- 308 [12] YOSINSKI, J., CLUNE, J., BENGIO, Y., AND LIPSON, H. How transferable are features in deep neural  
309 networks? In *Advances in Neural Information Processing Systems* (2014), Z. Ghahramani, M. Welling, C. Cortes,  
310 N. Lawrence, and K. Q. Weinberger, Eds., vol. 27, Curran Associates, Inc.
- 311 [13] ZHANG, Z., LUO, P., LOY, C. C., AND TANG, X. Facial landmark detection by deep multi-task learning. In  
312 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014), Springer, pp. 94–108.
- 313 [14] ZHAO, H., AND GORDON, G. Inherent tradeoffs in learning fair representations. In *Advances in Neural*  
314 *Information Processing Systems* (2019), H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and  
315 R. Garnett, Eds., vol. 32, Curran Associates, Inc.

# 316 Appendices

## 317 A Variational Autoencoder Architecture

Layer Name	Output Size	Settings
conv1	$16 \times 32 \times 32$	filter: $3 \times 3$ ; stride: 1; padding: 1
bn1	$16 \times 32 \times 32$	features: 16
ReLU	$16 \times 32 \times 32$	
conv2	$32 \times 16 \times 16$	filter: $3 \times 3$ ; stride: 2; padding: 1
bn2	$32 \times 16 \times 16$	features: 32
ReLU	$32 \times 16 \times 16$	
conv3	$32 \times 16 \times 16$	filter: $3 \times 3$ ; stride: 1; padding: 1
bn3	$32 \times 16 \times 16$	features: 32
ReLU	$32 \times 16 \times 16$	
conv4	$16 \times 8 \times 8$	filter: $3 \times 3$ ; stride: 2; padding: 1
bn4	$16 \times 8 \times 8$	features: 16
ReLU	$16 \times 8 \times 8$	
Flatten	1024	
fc1	512	
bn5 1D	512	features: 512
ReLU	512	
<b>mu</b>	64	
<b>logvar</b>	64	

Table 2: Encoder architecture. The convolutional and batchnorm layers are standard 2D versions unless otherwise indicated. The last two layers are the output layers for respectively the mean and log variance of the distribution. These both take the last ReLU layer as input. Batch size dimension is omitted.

Layer Name	Output Size	Layer
fc1	512	
ReLU	512	
fc2	1024	
bn1 1D	1024	features: 1024
ReLU	1024	
UnFlatten	$16 \times 8 \times 8$	
convT1	$32 \times 16 \times 16$	filter: $3 \times 3$ ; stride: 2; padding: 1; output padding: 1
bn2	$32 \times 16 \times 16$	features: 32
ReLU	$32 \times 16 \times 16$	
convT2	$32 \times 16 \times 16$	filter: $3 \times 3$ ; stride: 1; padding: 1; output padding: 0
bn3	$32 \times 16 \times 16$	features: 32
ReLU	$32 \times 16 \times 16$	
convT3	$16 \times 32 \times 32$	filter: $3 \times 3$ ; stride: 2; padding: 1; output padding: 1
bn4	$16 \times 32 \times 32$	features: 16
ReLU	$16 \times 32 \times 32$	
convT4	$1 \times 32 \times 32$	filter: $3 \times 3$ ; stride: 1; padding: 1; output padding: 0
bn5	$1 \times 32 \times 32$	features: 1
Sigmoid	$1 \times 32 \times 32$	

Table 3: Decoder architecture. Transposed Convolutions, indicated with *convT*, are used as a inverse of standard convolutional layers. The architecture was chosen to resemble the inverse of the encoder. The final non linear activation is a sigmoid since this matches the loss function and original image value range. Batch size dimension is omitted.

318 **B Discriminator Architecture**

Layer Name	Output Size	Layer
fc1	32	
ReLU	32	
Dropout	32	Keep prob: 0.5
fc2	32	
ReLU	32	
Dropout	32	Keep prob: 0.5
fc3	1	
Sigmoid	1	

Table 4: Discriminator architecture. The input are cluster assignment per sample of shape  $B \times C$ , where  $B$  is the batch size and  $C$  the number of clusters.

319 **C Hyperparameters**

Model	Parameter	Range	Optimal VAE	Optimal ResNet
<b>VAE</b>	batch size	64, 128, 256	64	-
	learning rate	1e-4, 1e-3, 1e-2, 1e-1	0.001	-
	epochs	20, 50, 100, 500	50	-
<b>DEC</b>	batch size	64, 128, 256	256	128
	iterations	10000 – 60000	50000	20000
	learning rate	1e-4, 1e-3, 1e-2, 1e-1, 1, 10	0.0001	0.0001
<b>DFC</b>	batch size	64, 128, 256	64	128
	iterations	20000, 35000, 50000	20000	20000
	learning rate	1e-4, 1e-3, 1e-2, 1e-1	0.001	0.001

Table 5: Complete list of hyperparameter search space and the optimal values. The search was performed with bayesian and a grid search for the datasets that rely on the custom VAE. Because of resource limitations the hyperparameter search for the more complex datasets that rely on ResNet, the search was manually and less extensive performed.

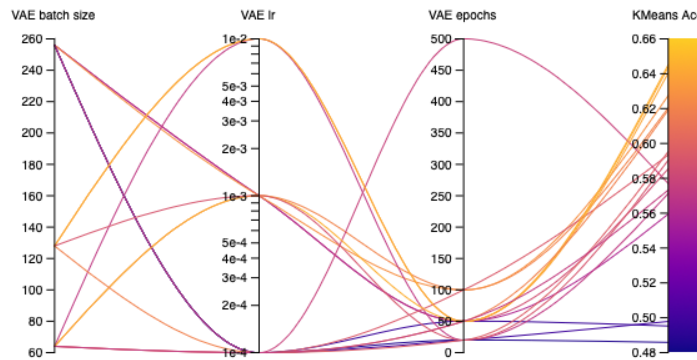


Figure 5: Overview of the results for the VAE hyperparameter search, evaluated with the clustering performance of a KMeans module on MNIST-USPS.

Model	Parameter	Value
<b>KMeans</b>	n init	20
	max steps	5000
	init	k-means++
	tol	1e-4
	precompute distances	auto
<b>TSNE</b>	algorithm	auto
	n components	2
	perplexity	30
	early exaggeration	12
	learning rate	200
	n iter	1000
	metric	euclidean
	init	random
	method	barnes hut
<b>VAE &amp; DEC</b>	angle	0.5
	Adam betas	(0.9, 0.999)
<b>DFC</b>	Adam weight decay	0
	Adam betas	(0.9, 0.999)
	Adam weight decay	1e-4

Table 6: Complete list of default hyperparameters per model not included in the hyperparameter search.

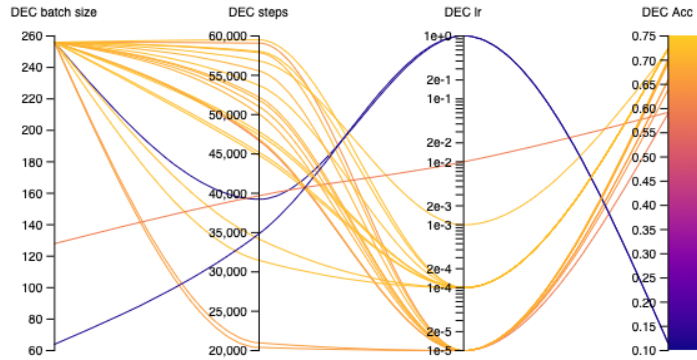


Figure 6: Overview of the results for the DEC hyperparameter search using the best performing pretrained VAE on MNIST-USPS.

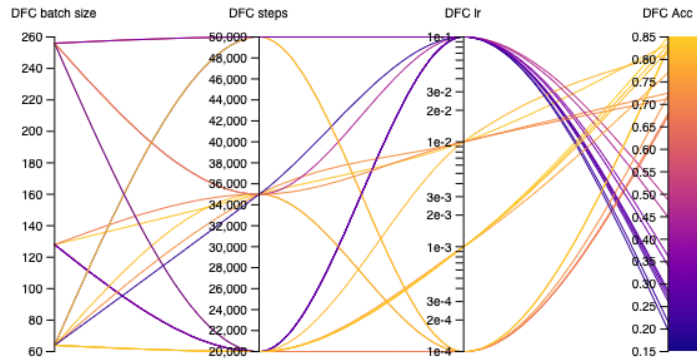


Figure 7: Overview of the results for the DFC hyperparameter search using the best performing pretrained VAE and DEC evaluated on MNIST-USPS.