

# Are X-ray landmark Detection Models fair? A preliminary assessment and mitigation strategy

Anonymous ICCV submission

Paper ID 10

## Abstract

001 *Datasets used for benchmarking are always acquired with*  
 002 *a view to representing different categories equally, with the*  
 003 *best intentions to be fair to all. Whilst it is usually as-*  
 004 *sumed that equal numerical representation in the training*  
 005 *data leads to similar accuracy among demographic groups,*  
 006 *so far, there has been next to no investigation or measure-*  
 007 *ment of this assumption for the anatomical landmark de-*  
 008 *tection task. In this work, we define what it means for*  
 009 *anatomical landmark detection to be carried out fairly on*  
 010 *different demographic categories, evaluating the fairness of*  
 011 *models trained on two publicly available X-ray datasets that*  
 012 *are known to be balanced, and showing how unfair predic-*  
 013 *tions can uncover metadata attributes intended to be hid-*  
 014 *den. We further design a potential mitigation strategy in the*  
 015 *landmark detection context, adapting a group optimization*  
 016 *method typically employed for debiasing image classifica-*  
 017 *tion models, obtaining a partial improvement in terms of*  
 018 *per-keypoint fairness, while paving the way for further re-*  
 019 *search in this field.*

## 020 1. Introduction

021 Precise and reliable anatomical landmark detection is criti-  
 022 cal for several clinical tasks [2, 9]. While the focus has been  
 023 on overall accuracy [4, 5, 10] and confidence [11], few stud-  
 024 ies have addressed bias within these models [3]. Biased and  
 025 unfair predictions in medical imaging can stem from non-  
 026 representative training datasets or from models that inadver-  
 027 tently perform better for certain demographic sub-groups  
 028 (i.e., age, gender, race). Limited literature exists on fairness  
 029 assessment in landmark detection, and is mainly for face  
 030 recognition datasets [8]. However, fairness in anatomical  
 031 landmark prediction remains largely unexplored, despite its  
 032 crucial clinical applications, such as diagnosis and surgical  
 033 treatment planning [15]. Our work addresses this critical  
 034 gap by establishing a protocol for assessing fairness in X-  
 035 ray anatomical landmark prediction. Our main contribution

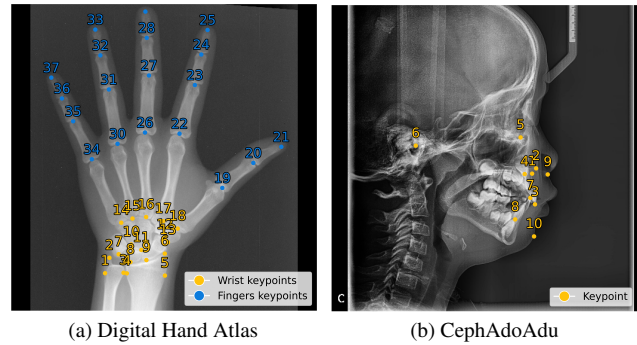


Figure 1. Numbered ground truth landmarks annotations for considered datasets.

is to show how fairness must be evaluated at *single keypoint*  
 level (see Fig. 1), since global measures hide fairness issues that may only affect a specific subset of keypoints. To this end, we adapt a popular classification fairness metric for use with landmark detection, further investigating the relationship between landmark prediction and patient metadata (age, gender), we show how errors on keypoints can be used to infer sensitive attributes, potentially raising privacy concerns. After measuring the fairness issue, we propose a potential mitigation approach based on a group optimization method typically employed for debiasing image classification models [13, 18], that is GroupDRO [16]. Our results show a partial improvement in the fairness metrics with negligible degradation of the overall landmark detection accuracy. To our knowledge, we are the first to expose a potential lack of fairness in the context of anatomical landmark detection, which occurs even when the training data is carefully acquired in balanced categories. This work is put forward as a critical foundation for improving data acquisition makeup and for developing benchmarking criteria. At the same time, we aim to shed light on the necessity of developing an ad-hoc solution for improving fairness in the context of anatomical landmark detection.

059 **2. Approach**

060 **2.1. Reference Datasets**

061 Since a study on attribute bias requires plentiful raw images and metadata, of the publicly available contenders (described in [5, 19]), only the Digital Hand Atlas (DHA) [7] and the CephAdoAdu dataset are fit for purpose.

062 The **DHA dataset** (Fig. 1a) includes 909 radiographs (average size:  $1563 \times 2169$  pixels) annotated with 37 landmarks. Among the available demographic attributes, we consider age and gender, which divide patients into groups large enough to allow a reliable assessment of group fairness. The dataset is balanced by design, with equal male and female patients per age group and broadly similar numbers of patients per age group. Ages range from 9 to 18 years, but due to small per-group sizes, we cluster them into younger’ (9–13 y.o.) and older’ (14–18 y.o.). The **CephAdoAdu dataset** (Fig. 1b) is a new benchmark comprising cephalometric X-ray images across age groups. Our dataset version includes 350 adult and 350 adolescents X-ray images, and manually annotated with 10 key landmarks. The training protocol is age-balanced, with 400 images (including 40 for validation) for training and 300 for testing, ensuring even distribution of adult and adolescent cases.

082 **2.2. Machine learning framework**

083 We address landmark detection by framing it as a supervised pixel-wise classification task to produce output heatmaps. Specifically, we exploit a U-Net model to predict several landmarks at once, creating one heatmap for each landmark as separate output channels ( $n$  heatmaps, where  $n$  is the number of annotated landmarks). As per [11], we generate ground-truth heatmaps from a single pixel input annotation. Formally,  $H_s(i, j) = \mathbb{1}(i = x_s \wedge j = y_s)$  where  $H_s(i, j)$  is the heatmap value at pixel  $(i, j)$ , the ground truth coordinates of the landmark ( $s$ ) are at pixel  $(x_s, y_s)$ , and  $\mathbb{1}$  evaluates to 1 only when the condition is satisfied. The output heatmap intensities are in  $[0, 1]$ , the hottest of which is the predicted landmark location.

096 **2.3. Adapting fairness evaluation metrics**

097 Since the dataset is designed to be balanced, it is assumed to be fair across all considered demographic categories (*gender* and *age*). Our goal is to identify the presence of group fairness issues [12] within the generic task of landmark detection. It is crucial to define an appropriate way of measuring whether a landmark detection model is *fair* (or *unfair*). A popular fairness metric for classification tasks is the *Demographic Parity* (DP), a.k.a. Statistical Parity [1, 6]. DP measures whether predicting a positive outcome is independent of a certain sensitive attribute. In a binary classification task, given a training dataset  $\mathcal{D} = \{(x_1, y_1, g_1), \dots, (x_n, y_n, g_n)\}$ , where

$y \in \{0, 1\}$  is the target label, and  $g$  encodes a group ( $0 \rightarrow$  Male,  $1 \rightarrow$  Female), DP is satisfied when a positive outcome is equal across different demographic groups:

$$P(y=1 | g=0) = P(y=1 | g=1) \quad (1)$$

In this setting, this can be verified by computing the classifier’s True Positive Rate  $TPR := \frac{TP}{TP+FN}$ , separately for each group. The largest absolute  $TPR$  difference between group pairs provides an empirical measure of a classifier’s fairness.

In landmark detection tasks, instead of being right or wrong, the prediction error is measured through the *Mean Radial Error* (MRE); this is the average distance between the predicted and the actual landmark positions, averaged over all landmarks. Here we use the Euclidean ( $L_2$ ) distance. Another measure of the model’s accuracy is the *Success Detection Rate* (SDR), reporting the proportion of predicted landmarks within a clinically acceptable distance threshold from the ground truth. In the case of  $n$  landmarks and a threshold  $\phi$ ,  $SDR = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(MRE_i \leq \phi)$ . Evidently, SDR can be translated as the TPR for a landmark prediction task: a True Positive is counted when a landmark is predicted within a threshold  $\phi$  from the ground truth. Otherwise, the model is said to have missed the prediction (False Negative). In our analysis, we compute the DP for each keypoint (KP), to uncover potential fairness issues hiding in individual keypoint predictions. This results in a specific DP value for each landmark, calculated with respect to all available sensitive attributes.

137 **2.4. Mitigating Fairness Issues**

In this work, we provide an attempt at mitigating the emerged fairness issues, which seem to be related not only to demographic groups but also stemming from particular anatomical keypoints. As such, we opt to address the problem with a subpopulation characterization as fine-grained as possible. Given a dataset with  $K$  keypoints and  $G$  known demographic groups, we consider  $K \times G$  possible subgroups. We encode them as an additional set of labels  $\mathcal{G} = \{0, \dots, G\}$ . For an input image  $x \in \mathbb{R}^{C \times W \times H}$  and per-keypoint ground-truth locations  $y \in \mathbb{R}^{K \times 2}$ , we provide a group label  $g \in \mathcal{G}^K$ , considering each keypoint separately, separating loss contributions from every  $G$  demographic group. For instance, in the case of the DHA dataset, we consider 37 keypoints and 4 demographic groups: {Young Males, Young Females, Old Males, Old Females}, for a total of 148 subgroups, where  $g = 0$  denotes KP1 from Young Males, and  $g = 148$  denotes KP37 in Old Females. Our subclass categorization allows us to frame the learning problem as

$$\hat{\theta} := \arg \min_{\theta \in \Theta} \max_{g \in \mathcal{G}} \left\{ \mathbb{E}_{(x,y) \sim \hat{P}_g} [\ell(\theta; (x, y))] \right\} \quad (2)$$

158 where  $\theta$  are the Unet parameters to be optimized. This ob-  
 159 jective, known as GroupDRO [16], is originally intended for  
 160 mitigating spurious correlations and improving worst-group  
 161 generalization in image classification settings. We provide  
 162 a customized adaptation of the method for landmark detec-  
 163 tion, fine-tuning the original model trained without any mit-  
 164 gation strategy with the objective in Equation 2.

### 165 3. Experiments

#### 166 3.1. Experiment details

167 Experiments utilize a U-Net with an ImageNet pre-trained  
 168 DenseNet121 encoder. Images are padded and resized to  
 169  $512 \times 512$  pixels maintaining aspect ratio and normalized to  
 170  $[0, 1]$ . Optimization employs AdamW for up to 200 epochs  
 171 with early stopping. The learning rate starts at  $10^{-3}$ , ad-  
 172 justed by Exponential scheduler. The batch size is 8 with  
 173 a gradient accumulation of 8. For SDR computation,  $\phi$  is  
 174 set to 2 mm, as it is the more restrictive threshold gener-  
 175 ally reported [5]. We compute metrics by converting pixel  
 176 distances to millimeters: for CephAdoAdu, we use a pixel  
 177 resolution of 0.1 mm; for DHA, we assume a 50 mm dis-  
 178 tance between wrist endpoints, as proposed by [14].

##### 179 3.1.1. Baseline computation

180 As a baseline, we perform 10 different hold-outs of the data,  
 181 preserving the balance between the demographic groups.  
 182 For each hold-out, we train a model on the training set  
 183 and evaluate it on the held-out test set. Finally, we report  
 184 the mean and standard deviation for each evaluation metric  
 185 across the ten runs. Fig. 2a (top) shows the average MRE  
 186 values for each keypoint of the DHA dataset across the  
 187 ten runs. Wrist keypoints (KP1 to KP18) generally exhibit  
 188 higher MRE values, indicating that they are somewhat more  
 189 challenging to detect, with several keypoints far exceeding  
 190 the overall average. Finger keypoints have comparatively  
 191 better performance, though some keypoints (KP19, KP36)  
 192 still exhibit high MRE. The overall MRE across all key-  
 193 points in the 10 performed runs is  $0.72\text{ mm}$ , with the MRE  
 194 for the wrist keypoints being slightly higher ( $0.84\text{ mm}$ )  
 195 compared to MRE for finger keypoints ( $0.61\text{ mm}$ ). Fig.  
 196 2b shows the same analysis replicated for the CephAdoAdu  
 197 dataset. In this dataset, the overall MRE is  $1.13\text{ mm}$ , with  
 198 some keypoints harder to detect than the average (e.g., KP4,  
 199 KP6 and KP8). For both datasets, the high variability in  
 200 MRE across keypoints propagates in the SDR computations  
 201 and highlights the importance of considering each keypoint  
 202 individually, as relying solely on metrics averaged across  
 203 the keypoints could easily hide any detection issue poten-  
 204 tially correlated to demographic groups.

##### 205 3.1.2. Fairness assessment

206 First, we compute the adapted fairness metric with respect  
 207 to specific sensitive attributes over the 10 hold-outs pre-

208 viously introduced and averaged across all the available  
 209 keypoints. For the DHA dataset, we get an overall maxi-  
 210 mum DP equal to  $0.045 \pm 0.009$ , while for the CheAdoAdu  
 211 dataset, an average value of  $0.080 \pm 0.006$  is obtained. Such  
 212 relatively low values would suggest that our models are fair  
 213 to the sensitive attributes. However, motivated by our pre-  
 214 vious results, whilst identifying peaks and troughs in the  
 215 MRE for specific keypoints, we deepen the fairness eval-  
 216 uation by framing our analysis as a per-keypoint problem.  
 217 Fig.2a (bottom) and Fig.2b (bottom) show the same analysis  
 218 on the sensitive attributes for single keypoints for our two  
 219 datasets. Analyzing these figures reveals significant vari-  
 220 ations in DP values across keypoints with respect to dif-  
 221 ferent sensitive attributes. Starting from the DHA dataset,  
 222 for example, KP1 and KP3 (ulna) show a DP value of 0.20,  
 223 meaning a maximum gap of 20% in the SDR across de-  
 224 mographic groups. Specifically, this maximum difference  
 225 arises between *female* patients in the two *age* groups. Inter-  
 226 estingly, this doesn't align with corresponding medical re-  
 227 search [17], who find no statistically significant difference  
 228 between ages in males and females, ultimately suggesting  
 229 dataset bias. Moreover, wrist keypoints show a DP value on  
 230 average much higher than fingers, suggesting higher fair-  
 231 ness in finger regions and underscoring the need for individ-  
 232 ual keypoint analysis over averaging. To further investigate  
 233 the statistical significance of the obtained results, for each  
 234 of our ten runs, we perform a metadata attribute random-  
 235 ization experiment. Specifically, we shuffle the sensitive  
 236 attributes with a probability of 50% for each sample, ob-  
 237 taining by construction a test evaluation uncorrelated with  
 238 metadata attributes. We report the corresponding average  
 239 DP per keypoint in orange in Figs. 2a and 2b. For the DHA  
 240 dataset, the most unfair wrist keypoints exhibit a DP sig-  
 241 nificantly higher than the attribute-randomized counterpart,  
 242 while finger keypoints, unaffected by fairness issues, main-  
 243 tain similar values across both settings. In the CephAdoAdu  
 244 dataset, KP4, KP5, and KP6 show higher demographic par-  
 245 ity but not significantly above the randomized counterpart.  
 246 Notably, KP1 has a DP of 0.17, exceeding the randomized  
 247 experiment, suggesting a potential fairness issue. Finally,  
 248 to further support the obtained results, excluding that identi-  
 249 fied fairness issues are a consequence of a specific or subop-  
 250 timal model, we report a comparison with available State-  
 251 Of-The-Art (SOTA) and perform an ablation study on the  
 252 Unet encoder in Table 1. Different models show similar val-  
 253 ues for the average DP across keypoints and roughly similar  
 254 average MREs, being competitive with SOTA.

##### 255 3.1.3. Results on fairness mitigation

256 Fig. 3a shows the results of our model fine-tuned with  
 257 the GroupDRO objective in terms of DP for the DHA  
 258 dataset (top) and CephAdoAdu (bottom). Regarding the  
 259 DHA dataset, the proposed mitigation approach brings a  
 260 general decreasing of DP across keypoints. However, the

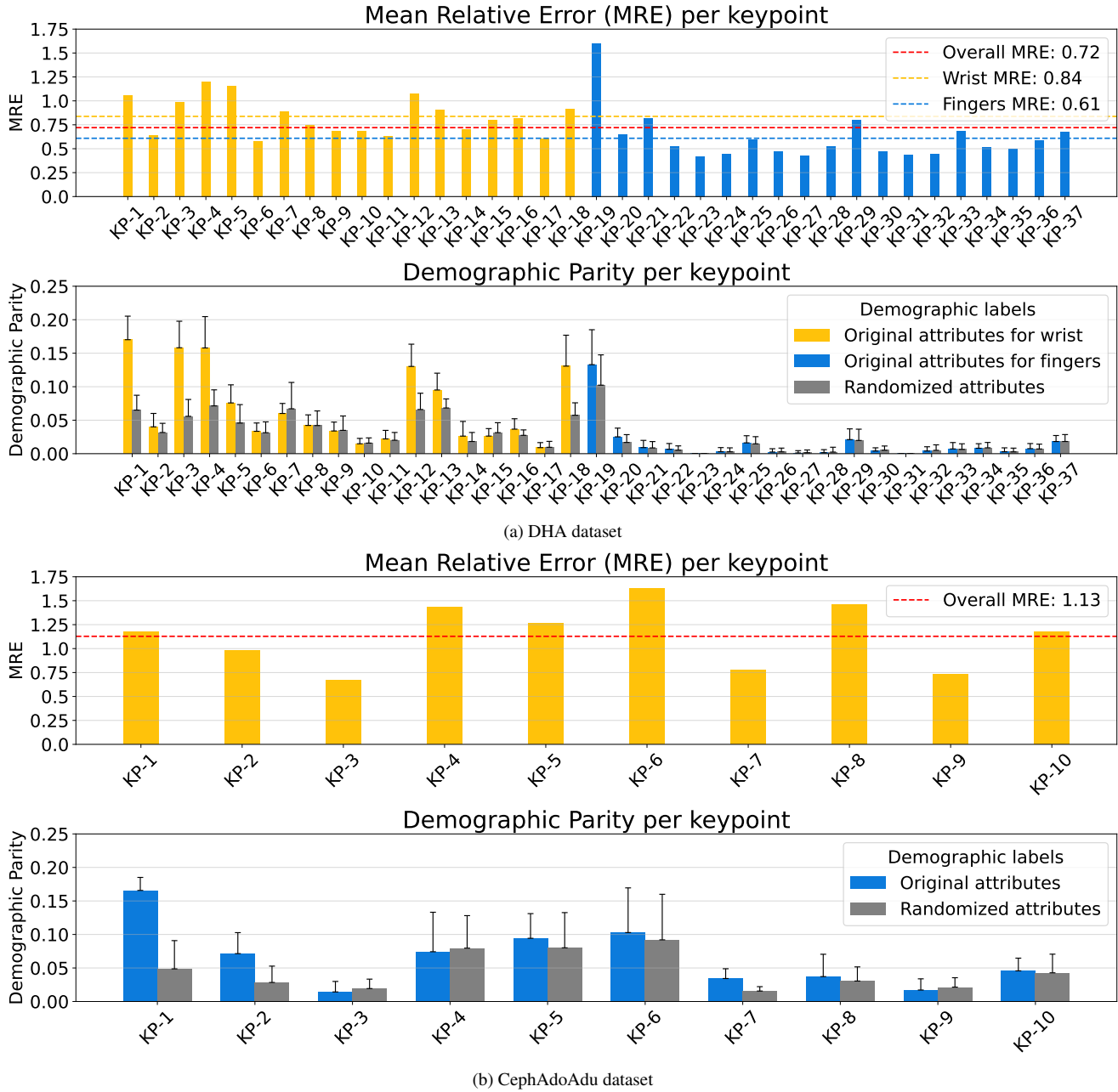


Figure 2. Assessment of keypoint prediction errors and demographic parity. (a) Top panel: MRE for 37 DHA dataset keypoints, with higher error in wrist vs. finger keypoints. Bottom panel: Demographic Parity measurements with original and randomized attributes across folds. (b) Top panel: MRE for 10 CephAdoAdu dataset keypoints. Bottom panel: Demographic Parity measurements.

261 fairness issue is not entirely solved, with some wrist key-  
 262 points yet presenting a final DP higher than 0.10 (e.g., KP1,  
 263 KP11, KP18 and KP19). A similar trend is observed in the  
 264 CephAdoAdu dataset, with some keypoints improving their  
 265 DP (e.g., KP4 and KP8). Again, the fairness issue is not  
 266 entirely solved, with KP1 still presenting a DP higher than

0.15. Importantly, the mitigated models roughly preserve  
 the average MRE across keypoints for both datasets, with a  
 maximum drop of 0.07 and 0.04 for the CephAdoAdu and  
 the HDA dataset, respectively.

267  
 268  
 269  
 270

Table 1. Top. Comparison of state-of-the-art results for anatomical landmark detection in X-ray images. Bottom. Ablation on specific Unet backbone.

| Methods            | CephAdoAdu          |          |       |       |       |                    | Digital Hand Atlas  |          |       |       |                          |                            |
|--------------------|---------------------|----------|-------|-------|-------|--------------------|---------------------|----------|-------|-------|--------------------------|----------------------------|
|                    | MRE ↓<br>(mm, std.) | SDR(%) ↑ |       |       |       | DP ↓<br>(avg. KPs) | MRE ↓<br>(mm, std.) | SDR(%) ↑ |       |       | DP Wrist ↓<br>(avg. KPs) | DP Fingers ↓<br>(avg. KPs) |
|                    |                     | 2mm      | 2.5mm | 3mm   | 4mm   |                    |                     | 2mm      | 4mm   | 10mm  |                          |                            |
| SCN [14]           | 1.73 (1.06)         | 82.97    | 90.40 | 93.37 | 96.57 | -                  | 0.66                | 94.99    | 99.27 | 99.99 | -                        | -                          |
| GU2Net [20]        | 1.69 (0.91)         | 80.33    | 88.13 | 91.47 | 95.57 | -                  | 0.84                | 95.40    | 99.35 | 99.75 | -                        | -                          |
| CeLDA [19]         | 1.05 (0.33)         | 89.13    | 93.60 | 96.17 | 98.67 | -                  | -                   | -        | -     | -     | -                        | -                          |
| Ours (resnet50)    | 1.13 (0.04)         | 85.90    | 91.43 | 94.43 | 97.50 | 0.062 (0.004)      | 0.80 (0.02)         | 96.47    | 99.16 | 99.69 | 0.076 (0.006)            | 0.034 (0.007)              |
| Ours (vgg19)       | 1.10 (0.01)         | 85.77    | 90.83 | 93.93 | 96.97 | 0.065 (0.003)      | 1.04 (0.20)         | 95.64    | 98.48 | 99.25 | 0.080 (0.004)            | 0.029 (0.002)              |
| Ours (densenet121) | 1.12 (0.04)         | 86.97    | 91.50 | 94.57 | 97.73 | 0.081 (0.006)      | 0.76 (0.06)         | 97.05    | 99.52 | 99.88 | 0.073 (0.011)            | 0.017 (0.008)              |

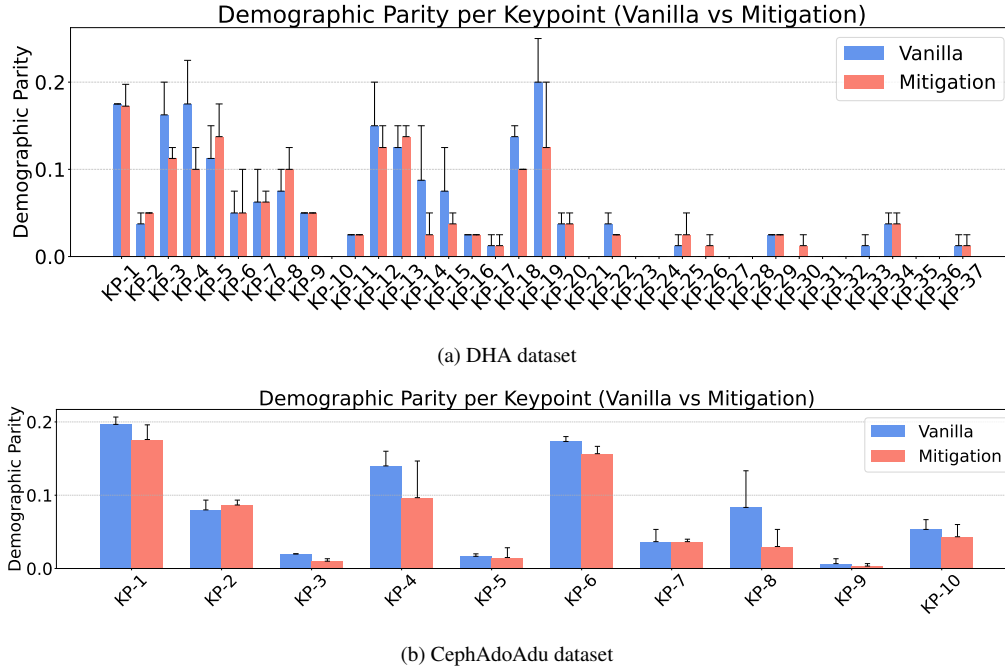


Figure 3. Per key-point Demographic Parity (DP) on vanilla and fairness mitigated models for the (a) DHA and (b) CephAdoAdu datasets.

Table 2. 5-fold classification accuracy of a CNN trained on images and an RF trained on MREs, for various attributes.

| Dataset    | Sensitive Attribute | Filtered Attribute | CNN image Classifier | RF MRE-based Classifier |
|------------|---------------------|--------------------|----------------------|-------------------------|
| DHA        | Age                 | male               | 0.53 ± 0.08          | 0.68 ± 0.05             |
|            |                     | female             | 0.56 ± 0.07          | 0.64 ± 0.12             |
|            | Gender              | young              | 0.56 ± 0.07          | 0.73 ± 0.13             |
|            |                     | old                | 0.55 ± 0.08          | 0.72 ± 0.07             |
| CephAdoAdu | Age                 | None               | 0.59 ± 0.16          | 0.64 ± 0.05             |

271 **3.1.4. Privacy-related issues**

272 Our results show a correlation between the MRE on specific  
 273 keypoints and metadata attributes. Here, we evaluate if such  
 274 an undesired correlation is strong enough to infer the sensi-  
 275 tive attribute from the computed errors, potentially lead-  
 276 ing to a privacy issue. Specifically, for the CephAdoAdu

dataset we consider the only available attribute (age). For  
 the DHA dataset, where we have two sensitive attributes  
 (gender and age), we further filter data according to a spe-  
 cific metadata attribute (*young/old* and *female/male* respec-  
 tively), reported as *Filtered attribute* in Table 2. This ap-  
 proach prevents attribute mixing, isolating each sensitive  
 attribute’s contribution. Thus, we train a Random Forest  
 (RF) classifier, exploiting the MREs corresponding to each  
 keypoint as features and the target sensitive attribute as la-  
 bels. We perform a 5-fold cross-validation, replicating the  
 same experiments considering the test MRE across all key-  
 points as input features. Table 2 summarizes the obtained  
 results. For both datasets, the MREs across keypoints bring  
 an average test accuracy much higher than a random guess,  
 with a maximum value of 0.75 for the sensitive attribute *age*  
 in the *female* filtered attribute for the HDA and 0.64 for the

293 CephAdoAdu dataset. To ensure that these results are not a  
294 simple consequence of the sensitive attribute being inferred  
295 from the images, we train a CNN directly on the X-ray im-  
296 ages with the same folds. As we can see in Table 2, we  
297 obtain an accuracy close to random guessing, further prov-  
298 ing that the results are an actual consequence of the fairness  
299 issue.

#### 300 4. Conclusions and future work

301 Despite the best intentions to acquire and anonymise patient  
302 data, we uncover concerns around the varying performance  
303 of landmark detection models known to be performing well.  
304 Privacy can be compromised through unintentional lack of  
305 fairness in such model–data pairs. Further work is required  
306 to understand this phenomenon better, potentially requir-  
307 ing the acquisition of new datasets and experimenting with  
308 different proportions of subjects in each demographic cate-  
309 gory, with a view to stabilising demographic parity. In this  
310 work, we adapt a typical mitigation strategy for image clas-  
311 sification model debiasing, obtaining a partial mitigation of  
312 the described phenomenon. Despite promising results, our  
313 work eventually aims to highlight the necessity of design-  
314 ing ad-hoc methods (e.g., involving domain practitioners to  
315 define proper anatomical priors) for mitigating unfairness in  
316 anatomical landmark detection, potentially paving the way  
317 for multiple future investigations.

#### 318 References

319 [1] Richard J Chen, Judy J Wang, Drew FK Williamson,  
320 Tiffany Y Chen, Jana Lipkova, Ming Y Lu, Sharifa Sahai,  
321 and Faisal Mahmood. Algorithmic fairness in artificial in-  
322 telligence for medicine and healthcare. *Nature biomedical*  
323 *engineering*, 7(6):719–742, 2023. 2

324 [2] Allison Clement, Abhinav Singh, and Irina Voiculescu.  
325 Landmark-based screening: Femoral head coverage and graf  
326 classification in infant developmental dysplasia of the hip. In  
327 *European Conference on Computer Vision (ECCV)*. Woman  
328 in Computer Vision (WiCV) Workshop, Springer Cham,  
329 2024. 1

330 [3] Li David, Lin Cheng Ting, Sulam Jeremias, and Yi Paul H.  
331 Deep learning prediction of sex on chest radiographs: a po-  
332 tential contributor to biased algorithms. *Emergency Radi-*  
333 *ology*, 29(2):365–370, 2022. © 2022. American Society of  
334 Emergency Radiology. 1

335 [4] Roberto Di Via, Francesca Odone, and Vito Paolo Pas-  
336 tore. Self-supervised pre-training with diffusion model  
337 for few-shot landmark detection in x-ray images. *CoRR*,  
338 abs/2407.18125, 2024. 1

339 [5] Roberto Di Via, Matteo Santacesaria, Francesca Odone, and  
340 Vito Paolo Pastore. Is in-domain data beneficial in trans-  
341 fer learning for landmarks detection in x-ray images? In  
342 *IEEE International Symposium on Biomedical Imaging, ISBI*  
343 *2024, Athens, Greece, May 27-30, 2024*, pages 1–5. IEEE,  
344 2024. 1, 2, 3

[6] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Rein- 345  
gold, and Richard Zemel. Fairness through awareness. In 346  
*Proceedings of the 3rd innovations in theoretical computer*  
*science conference*, pages 214–226, 2012. 2 347  
348

[7] Arkadiusz Gertych, Aifeng Zhang, James W. Sayre, Sylwia 349  
Pospiech-Kurkowska, and H. K. Huang. Bone age assess- 350  
ment of children using a digital hand atlas. *Comput. Medical*  
*Imaging Graph.*, 31(4-5):322–331, 2007. 2 351  
352

[8] Leonardo Iurada, Silvia Bucci, Timothy M. Hospedales, 353  
and Tatiana Tommasi. Fairness meets cross-domain learn- 354  
ing: a new perspective on models and metrics. *CoRR*, 355  
abs/2303.14411, 2023. 1 356

[9] Yankun Lang, Xiaoyang Chen, Hannah H. Deng, Tianshu 357  
Kuang, Joshua C. Barber, Jaime Gateno, Pew-Thian Yap, 358  
and James J. Xia. Dentalpointnet: Landmark localization on 359  
high-resolution 3d digital dental models. In *Medical Image*  
*Computing and Computer Assisted Intervention - MICCAI*  
*2022 - 25th International Conference, Singapore, September*  
*18-22, 2022, Proceedings, Part II*, pages 444–452. Springer,  
2022. 1 360  
361  
362  
363  
364

[10] Juneja Mamta, Garg Poojita, Kaur Ravinder, Manocha Palak, 365  
Prateek, Batra Shivam, Singh Pradeep, Singh Shaswat, and 366  
Jindal Prashant. A review on cephalometric landmark detec- 367  
tion techniques. *Biomedical Signal Processing and Control*,  
66:102486, 2021. 1 368  
369

[11] James McCouat and Irina Voiculescu. Contour-hugging 370  
heatmaps for landmark detection. In *Proceedings of the*  
*IEEE/CVF Conference on Computer Vision and Pattern*  
*Recognition*, pages 20597–20605, 2022. 1, 2 371  
372  
373

[12] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina 374  
Lerman, and Aram Galstyan. A survey on bias and fairness in 375  
machine learning. *ACM computing surveys (CSUR)*, 54 376  
(6):1–35, 2021. 2 377

[13] Vito Paolo Pastore, Massimiliano Ciranni, Davide Marinelli, 378  
Francesca Odone, and Vittorio Murino. Looking at model 379  
debiasing through the lens of anomaly detection. In *Proceed-*  
*ings of the Winter Conference on Applications of Computer*  
*Vision (WACV)*, pages 2548–2557, 2025. 1 380  
381  
382

[14] Christian Payer, Darko Stern, Horst Bischof, and Martin 383  
Urschler. Integrating spatial configuration into heatmap re- 384  
gression based cnns for landmark localization. *Medical Im-*  
*age Anal.*, 54:207–219, 2019. 3, 5 385  
386

[15] WR Proffit, HW Fields, and DM Sarver. Contemporary or- 387  
thodontics: Elsevier health sciences. *Philadelphia, USA*,  
2006. 1 388  
389

[16] Shiori Sagawa\*, Pang Wei Koh\*, Tatsunori B. Hashimoto,  
and Percy Liang. Distributionally robust neural networks.  
In *International Conference on Learning Representations*,  
2020. 1, 3 390  
391  
392  
393

[17] E Sayit, A Tanrivermis Sayit, M Bagir, and Yüksel Terzi.  
Ulnar variance according to gender and side during aging:  
An analysis of 600 wrists. *Orthopaedics & Traumatology:*  
*Surgery & Research*, 104(6):865–869, 2018. 3 394  
395  
396  
397

[18] Nimit Sohoni, Jared Dunmon, Geoffrey Angus, Albert  
Gu, and Christopher Ré. No subclass left behind: Fine-  
grained robustness in coarse-grained classification problems.  
*Advances in Neural Information Processing Systems*, 33:  
19339–19352, 2020. 1 398  
399  
400  
401  
402

- 403 [19] Han Wu, Chong Wang, Lanzhuju Mei, Tong Yang, Min Zhu,  
404 Dinggang Shen, and Zhiming Cui. Cephalometric landmark  
405 detection across ages with prototypical network. In *Medi-*  
406 *cal Image Computing and Computer Assisted Intervention -*  
407 *MICCAI 2024 - 27th International Conference, Marrakesh,*  
408 *Morocco, October 6-10, 2024, Proceedings, Part V*, pages  
409 155–165. Springer, 2024. [2](#), [5](#)
- 410 [20] Heqin Zhu, Qingsong Yao, Li Xiao, and S. Kevin Zhou. You  
411 only learn once: Universal anatomical landmark detection.  
412 In *Medical Image Computing and Computer Assisted In-*  
413 *tervention - MICCAI 2021 - 24th International Conference,*  
414 *Strasbourg, France, September 27 - October 1, 2021, Pro-*  
415 *ceedings, Part V*, pages 85–95. Springer, 2021. [5](#)