# FairFace Challenge at ECCV 2020: Analyzing Bias in Face Recognition

Tomáš Sixta[1,*], Julio C. S. Jacques Junior[2,3,*], Pau Buch-Cardona[3,4], Neil M. Robertson[5], Eduard Vazquez[6], and Sergio Escalera[3,4]

[1] Czech Technical University in Prague, Czech Republic – `tomas.sixta@gmail.com`
[2] Universitat Oberta de Catalunya, Spain – `jsilveira@uoc.edu`
[3] Computer Vision Center, Spain
[4] Universitat de Barcelona, Spain – `sergio@maia.ub.es`
[5] The Queen's University of Belfast, United Kingdom – `N.Robertson@qub.ac.uk`
[6] Anyvision, United Kingdom – `eduardov@anyvision.co`

**Abstract.** This work summarizes the 2020 ChaLearn Looking at People Fair Face Recognition and Analysis Challenge and provides a description of the top-winning solutions and analysis of the results. The aim of the challenge was to evaluate accuracy and bias in gender and skin colour of submitted algorithms on the task of 1:1 face verification in the presence of other confounding attributes. Participants were evaluated using an in-the-wild dataset based on reannotated IJB-C, further enriched by 12.5K new images and additional labels. The dataset is not balanced, which simulates a real world scenario where AI-based models supposed to present fair outcomes are trained and evaluated on imbalanced data. The challenge attracted 151 participants, who made more than 1.8K submissions in total. The final phase of the challenge attracted 36 active teams out of which 10 exceeded 0.999 AUC-ROC while achieving very low scores in the proposed bias metrics. Common strategies by the participants were face pre-processing, homogenization of data distributions, the use of bias aware loss functions and ensemble models. The analysis of top-10 teams shows higher false positive rates (and lower false negative rates) for females with dark skin tone as well as the potential of eyeglasses and young age to increase the false positive rates too.

**Keywords:** face verification; face recognition; fairness; bias.

## 1 Introduction

Automatic face recognition is a general topic that includes both face identification and verification [29]. Face identification is the process of identifying someone's identity given a face image. This process is generally known as 1-to-n matching and could be seen as asking to the system "who is this person?". Face verification, on the other hand, is concerned with validating a claimed identity

---

* These (corresponding ✉) authors contributed equally to this work.

based on the image of a face, and either accepting or rejecting the identity claim (1-to-1 matching). A simple example of face verification is when people unlock their smartphones using their faces (e.g., authentication), whereas searching for the identity of a given individual in a database of missing people, for instance, could be an example of face identification.

Fairness in face recognition recently started to receive increasing interest from different segments of scientific communities [19,36,38,44]. This is partially due to the huge impact new technologies have in our daily lives. Face recognition has been routinely utilized by both private and governmental organizations around the world [16,53]. Automatic face recognition can be used for legitimate and beneficial purposes (e.g. to improve security) but at the same time its power and ubiquity heightens a potential negative impact unfair methods can have for the society [52,55,54,46]. Recently, these concerns led several major companies to suspend distribution of their products to US police departments until a legislation regulating its deployment is passed by US Congress [59,56,57].

Although not sufficient, a necessary condition for a legitimate deployment of face recognition algorithms is equal accuracy for all demographic groups. A gold standard for testing commercial products is the Face Recognition Vendor Test (FRVT) performed by National Institute of Standards and Technology (NIST) [23,5]. However, this test is not designed for iterative and fast evaluation of new research directions. There is also a growing number of works that evaluate the algorithms on public data [3,64,50,13] and are therefore limited by what data is available, i.e., typically either small scale high quality datasets or large scale datasets with noisy annotations.

To motivate research on fair face recognition and provide a new challenging accurately annotated dataset, we designed and ran a computational face recognition challenge where participants were asked to provide solutions that maximize both accuracy and two fairness scores (minimize bias score). The submissions were evaluated on a reannotated version of IJB-C [37] database, enriched by newly collected 12,549 public domain images. The dataset contains large variations in head pose, face size and other attributes (detailed in Sec. 4.1). The dataset is not balanced with respect to different attributes, which imposes another challenge for the participants and is intended to stimulate usage of bias mitigation methods, also because the final ranking is defined by a weighted combination of accuracy and fairness (giving the bias scores a higher weight). To this end, we propose a new evaluation metric derived from a causal model by means of a causal effect of protected attributes to the accuracy of the algorithm, detailed in Sec. 4.2. The challenge attracted a total of 151 participants, who made more than 1.8K submissions in total[1]. We expect the provided dataset and proposed fairness measure template to be a reference evaluation benchmark for face recognition systems, and that the outcomes of this challenge will help both to define priorities for future research as well as to help on the definition of technical requirements for real applications.

---

[1] Data and winning solutions codes are available at http://chalearnlap.cvc.uab.es/challenge/38/description

## 2   Ethics in Face Recognition

Face recognition methods have been researched for decades due to their wide number of scenarios for good[2]. They can be applied, e.g., in robotics, human-computer interaction, access and control, security, among others. Recently, face recognition research received additional attention due to the improved performance provided by deep learning architectures [24]. When it comes to public safety, past works raised the question about the efficacy of facial recognition systems for law enforcement following the apparent failure of the systems to identify suspects, reporting as possible reasons for failure problems like occlusions, angled facial shots, poor lighting or obscured facial features by hats or sunglasses [1]. However, recent studies show that automated methods for face analysis can also discriminate based on classes like gender and ethnicity [11], among others, which raised an additional focus of attention around such technologies. If face recognition methods are used to support decisions, erroneous but confident mis-identification can have serious consequences, and these possible and negative outcomes are making the society to rethink about what should be the limits of such technology, especially when it is applied at larger scales involving additional privacy concerns.

From a research point of view, a bottleneck to be solved is to develop methods that can work accurately for all target populations. While there is a need to promote good practices and reinforce regulations, we need to find the way to provide the required good (and fair) performance in practice, and if face recognition is to be applied, it should deal with the bias problem. Evidences show that the computer vision and machine learning research communities are starting to give visibility to different types of bias [10,22] and proposing different solutions to mitigate them (e.g., [58,11,21,6,65,68]). Nonetheless, additional efforts should be made to further reduce bias in future methods. This is precisely the main goal of the 2020 ChaLearn Looking at People Fair Face Recognition and Analysis Challenge, i.e., to stimulate and promote research on face recognition methods that produce fair outcomes.

## 3   Related Work

It is known that popular face recognition datasets like Labeled Faces in the Wild (LFW) [34], MegaFace [30], IJB-C [37], IMDB-WIKI [48,49], VGGFace2 [12] or MS-Celeb-1M [25] are imbalanced both in gender and skin colour [39]. To encourage research in fair face recognition there is growing number of datasets specifically designed with balance in mind and annotated for gender, ethnicity and potentially other attributes. Examples are Racial Faces in the Wild (RFW) [64] (40K images, 12K identities, subset of MS-Celeb-1M), Balanced Faces in

---

[2] For more information about ethics in AI you can visit the European guideline in the following link https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai.

the Wild (BFW) [47] (20K images, 0.8K identities, subjects sampled from VG-GFace2) or DiveFace [40] (150K images, 24K identities, subset of Megaface). Even though these datasets are important step towards fairer face recognition, using labels for ethnicity does not in general allow for comparing models across datasets, because unlike for skin colour [9] there is no widely accepted definition of ethnicity groups and the labels instead rely on judgment of the annotators. Furthermore, balancing alone may not be enough to guarantee fair models [2], which motivates research of bias mitigation methods.

Nowadays, the gold standard for evaluating accuracy and bias of face recognition algorithms is the ongoing FRVT Test performed by NIST [23,5]. The submitted (mostly commercial) algorithms are evaluated on four datasets composed of photographs from various visa/benefits US governmental applications. In total, there are 18.27 million images of 8.49 million people. Besides FRVT, there are numerous small scale evaluations of bias in publicly and commercially available algorithms (e.g. [64,50,13]) as well as analysis of bias in models trained from scratch on publicly available datasets [3], that in most cases report better accuracy for men and people with light skin colour.

Traditional measures of fairness are based on calculating certain statistics related to the error rate of the algorithm. For example, Equalized Odds requires the true positive and false positive rates to be equal for all protected groups (see [60,19] for a comprehensive review). These measures are easy to calculate, but without having background in statistics it can be difficult to choose the "correct" one for the task at hand. This is a serious shortcoming, because certain traditional measures in general contradict each other [8,15,31] as it was dramatically demonstrated on the case of COMPAS (a system used in some US states to predict the risk of recidivism) [45,18]. Individual Fairness [20] tries to overcome these shortcomings by deriving a fairness measure from intuition, that "*similar individuals are treated similarly*". However, it does not provide any general definition of similarity and only postpones the problem by proposing that it should be given by a regulatory body or a civil rights organization.

A growing number of state of the art measures is based on causal inference. They require a "model of the world" given as a causal diagram and the actual measure is then derived using this diagram, e.g. in terms of the causal effect of the protected attributes to the algorithm accuracy [69,41] or using counterfactuals [33], i.e., would the decision remain had the value of the protected attribute be different but everything else stayed the same. A crucial advantage of these approaches is that the underlying ethical views are encoded by the diagram in an easy to understand way, which exposes them to criticism and allows them to be changed if they prove to be inadequate. Furthermore, as these approaches are trying to identify the true causes of the unfairness, they can be used as a starting point for mitigating the bias in the real world.

Bias mitigation methods can be broadly divided based on what area of model deployment they target to pre-processing, in-processing and post-processing [7,43]. The most popular pre-processing technique is rebalancing the dataset [27,66], alternatively using synthetic data [32]. In-processing approaches include cost-

sensitive training (higher weights for underrepresented groups) [27], adversarial learning for removing the sensitive information from the features [4,66], tuning parameters of a loss function for different protected groups [40,63] or attempts to learn bias free representations in unsupervised way [61]. Examples of post-processing techniques are renormalizing the similarity score of two feature vectors based on the demographic groups of the corresponding images [51] or attaching more fully connected layers to the feature extractor in order to remove the sensitive information from the representations [40]. The FairFace Recognition challenge, described in Sec. 4, did not impose any constrain to the participants to what model stage bias mitigation should be addressed. The best solutions rely on a combination of different strategies, detailed in Sec. 5.

## 4   Challenge Design

The participants were asked to develop their face verification methods aiming for a reduced bias in terms of the protected attributes (i.e., gender and skin color). Developed methods needed to output a list of confidence scores given test ID pairs to be verified (higher score means higher confidence, that the image pair contains the same person). The challenge[3] was managed using Codalab[4], an open source framework for running competitions that allows result or code submission.

The challenge ran from 4th April to 1st July 2020, and included two different phases: development and test. In the development phase, the participants were provided with public train data (with labels) and validation data (without labels, from which they should make predictions). At the test stage, the validation labels were released to all participants as well as the test data (without labels, considered for the final evaluation). The challenge attracted a total of 151 registered participants. During development phase we received 1330 submissions from 48 teams, and 476 submission from 36 teams at the test stage, resulting in more than 1800 submissions in total. Additional schedule details and participation statistics are provided in the supplementary material.

### 4.1   The Dataset

The dataset used in the challenge is a reannotated version of IJB-C [37], further enriched by newly collected 12,549 public domain images. In total, there are 152,917 images from 6,139 identities. The images were annotated by Anyvision's internal annotation team for two protected attributes: gender (male, female) and skin colour (light corresponding to Fitzpatrick types I-III, dark corresponding to types IV-VI) and five legitimate attributes: age group (0-34, 35-64, 65+), head pose (frontal, other), image source (still image, video frame), wearing glasses and a bounding box size[5]. Detailed annotation instructions are in the supplementary

---

[3] https://competitions.codalab.org/competitions/24184

[4] https://competitions.codalab.org

[5] Attribute categories used in this work are imperfect for many reasons. For example, it is unclear how many skin colour and gender categories should be stipulated (or whether they should be treated as discrete categories at all). We base our definitions

material. Every attribute was annotated by at least 3 annotators (age and skin colour by 6, due to their subjectiveness, aiming to maximize the level of agreement). Labels for gender and skin colour were synchronized for each *identity* to the most prevalent ones, and labels of the other attributes were obtained by choosing for each *image* the most common label from the annotators.

For the purpose of the challenge, the dataset was split into training, validation and testing subsets containing 70%, 10% and 20% of identities. To facilitate evaluation of the submitted results we generated roughly half a million positive face image pairs (same identity) and half a million negative pairs for both validation and testing subsets. The pairs were selected such that the number of combinations of legitimate attributes is maximized. In the validation pairs there were 219 (positive) and 574 (negative) combinations, and test pairs contained 397 (positive) and 1162 (negative) combinations. Basic dataset statistics are summarized in Table 1. Few image samples and and statistics of the attributes are shown in Fig. 1 and Fig. 2, respectively.

**Table 1.** Dataset statistics.

|  | Train | Validation | Test | Total |
|---|---|---|---|---|
| **Images** | 100,186 | 17,138 | 35,593 | 152,917 |
| **Unique identities** | 4,297 | 614 | 1,228 | 6,139 |
| **Positive pairs** | - | 448,119 | 500,176 | - |
| **Negative pairs** | - | 552,672 | 500,963 | - |



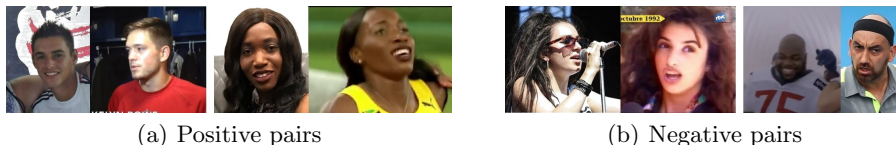(a) Positive pairs                    (b) Negative pairs

**Fig. 1.** Positive and negative samples of image pairs used in the challenge.

The images in the dataset have large variance in head pose, bounding box size and other attributes, which makes it challenging for face recognition. At the same time the distribution of these attributes is imbalanced, for example as seen in Fig. 2(a), there is considerably more white males that dark females. Such imbalances are common in real world datasets and we intentionally have not reballanced the data to encourage research of bias mitigation methods.

### 4.2   Evaluation Protocol

The challenge submissions were evaluated for bias in positive and negative pairs, and overall accuracy (given by AUC-ROC). The measure of bias/fairness was

---

on widely accepted traditional categories and our methodology and findings are expected to be applied later to any re-defined and/or extended attribute category.
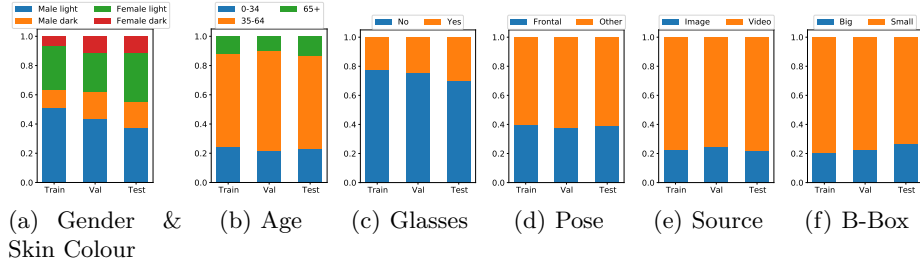
(a) Gender & Skin Colour    (b) Age    (c) Glasses    (d) Pose    (e) Source    (f) B-Box

**Fig. 2.** Distribution in percentage of attributes in training, validation and testing subsets of the dataset. Bounding box of a face was considered small if either its width or height was smaller than 224 px.

derived from a causal diagram shown in Fig. 3, in terms of a causal effect of protected attributes $A$ (gender and skin colour) to the output $\hat{Y}$ of the algorithm. The diagram was chosen using the following principle: the accuracy of the algorithm might be influenced (caused) directly by gender and skin colour but in addition there might be other variables that influence the accuracy and depend on the protected attributes. Some of these additional variables are seen as legitimate causes for different accuracy, whereas the others are proxies for unfair discrimination. It should be emphasized that the structure of the diagram and designations of the additional variables are not learned from the data but selected to express ethical views on the real world. This does not allow to select an objectively best diagram but instead provides transparency needed for the public to review it and potentially change it based on democratic discussion.
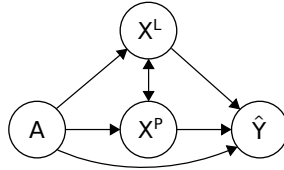


**Fig. 3.** Causal model used for our definition of fairness. $A$: protected attributes (gender and skin colour), $X^L$: legitimate attributes, $X^P$: proxy attributes, $\hat{Y}$: outcome of the algorithm. Note that in this challenge we deemed all additional attributes as legitimate, so $X^P$ did not contain any variables.

Our definition of fairness is inspired by intuition, that an algorithm is fair if given fixed values of the *legitimate* attributes its outcome remained the same regardless of the values of the protected and proxy attributes. Distinguishing legitimate and proxy attributes is crucial for the definition of fairness. By denoting an attribute legitimate we chose to ignore, that it might have different prevalence in different protected groups (i.e., break the causal link from the protected at-

tribute) and consequently that the algorithm has different accuracy for different groups. An example could be eyeglasses - as they can be easily removed, different accuracy caused by them is not seen as unfair even if they were worn more frequently by certain protected group. This is however not true for the proxy attributes - they are seen as mediators of potential unfair discrimination and therefore causal paths going over them must be included in the final objective.

Following the notation from [42], breaking causal links can be expressed by the $do()$ operator, which denotes an intervention on a variable. A prediction $\hat{Y}$ is fair with respect to protected attributes $A$ and causal diagram $M$ if for every pair of protected groups $a, a'$ and value $x^L$ of legitimate attributes

$$\sum_{X^P} P_M(\hat{Y}, X^P \mid do(A = a), do(X^L = x^L)) =$$
$$\sum_{X^P} P_M(\hat{Y}, X^P \mid do(A = a'), do(X^L = x^L)), \tag{1}$$

where $X^P$ denotes the proxy variables. As described in [42], $do(X)$ in diagram $M$ is equivalent to conditioning on plain $X$ in mutilated diagram $M^*$, where links leading to $X$ are removed. Furthermore, because in this challenge we deemed all additional attributes as legitimate, the criterion can be simplified to

$$P_{M^*}(\hat{Y} \mid A = a, X^L = x^L) = P_{M^*}(\hat{Y} \mid A = a', X^L = x^L). \tag{2}$$

As probability of error depends on a recognition threshold, we replace it by AUC-ROC metrics and use $AUC(a; x^L)$ to denote accuracy for positive pairs from protected group $a$ with legitimate attributes $x^L$ (all negative pairs are used as the negative samples for the ROC curve; accuracy for negative pairs is obtained in the same way with the roles of positive and negative samples reversed). To obtain a single numerical measure of bias, we define the discrimination $d(a; x^L)$ for protected group $a$ with legitimate attributes $x^L$ as a difference in accuracy for this group and the best one:

$$d(a; x^L) = \max_{a'} AUC(a'; x^L) - AUC(a; x^L). \tag{3}$$

The final measure of bias reported in the rankings is the difference between the average discriminations of the most and the least discriminated group:

$$Bias = max_a \frac{1}{|X|} \sum_{x^L} d(a; x^L) - min_a \frac{1}{|X|} \sum_{x^L} d(a; x^L). \tag{4}$$

### 4.3   Ranking Strategy

Having the accuracy and the two bias scores (for positive and negative pairs), participants were ranked by the average rank position obtained on each of these 3 variables. This way, bias is receiving more weight than accuracy. However, to prevent a random number generator from winning the competition we require that the accuracy of the submissions must be higher than the accuracy of our baseline model (see Sec. 4.4). Similarly, the submission of constant values would return Bias score = 0, due to the "$max - min$" strategy defined in Sec. 4.2.

### 4.4   The Baseline

We provide a baseline in order to set a reference point. We implemented a well-known standard solution for the face verification task based on a Siamese network [14] over a ResNet50 [26] backbone architecture (pretrained on Faces [12] database). Standard bounding box regression network for face detection was applied to detect the face region in every single image. Training pairs were generated by considering a subset of the dataset with highest possible diversity in terms of legitimate attributes. These pairs were fed to the model in balanced batches of 16 samples. The system was optimized with respect to maximizing only face verification accuracy confidence. As training strategy, only the layers from the 4th convolutional block of ResNet50 have been fine-tuned, using Adam as optimizer, $lr = 0.0001$ and Binary Cross-Entropy Loss, for 300 epochs.

## 5   Challenge Results, Winning Methods and Bias Analysis

### 5.1   The Leaderboard

Results obtained by the top-10 winning solutions at the development phase[6] are reported in Table 2. As it can be seen, results are very good if only accuracy is considered. Thus, the Bias scores can be considered a relevant tiebreaker factor, as one of the goals of the challenge is to stimulate research and development of fair face recognition methods.

**Table 2.** Top-10 solutions on the development phase (and Baseline results). The number inside the parenthesis indicate the global rank position for that particular variable, used to compute the average ranking.

| Participant | Average Ranking | Entries | Bias (+ pairs) | Bias (- pairs) | Accuracy |
|---|---|---|---|---|---|
| ustc-nelslip | 2.333333 (1) | 30 | 0.000142 (1) | 0.002956 (3) | 0.999287 (3) |
| zheng.zhu | 3.666667 (2) | 133 | 0.000344 (3) | 0.003781 (7) | 0.999442 (1) |
| CdtQin | 3.666667 (2) | 72 | 0.000472 (5) | 0.002334 (1) | 0.998477 (5) |
| crisp | 4.666667 (3) | 14 | 0.000935 (8) | 0.003193 (4) | 0.999394 (2) |
| haoxl | 4.666667 (3) | 73 | 0.000348 (4) | 0.003678 (6) | 0.998699 (4) |
| cam_vision | 5.000000 (4) | 95 | 0.000731 (6) | 0.002488 (2) | 0.995621 (7) |
| Hyg | 6.000000 (5) | 33 | 0.000814 (7) | 0.003305 (5) | 0.998402 (6) |
| senlin11 | 9.333333 (6) | 50 | 0.000165 (2) | 0.010091 (16) | 0.992093 (10) |
| hanamichi | 10.666667 (7) | 91 | 0.001631 (9) | 0.006760 (10) | 0.987382 (13) |
| paranoidai | 12.000000 (8) | 156 | 0.003779 (12) | 0.007745 (13) | 0.988359 (11) |
| *Baseline* | 38.333333 (33) | 1 | 0.057620 (40) | 0.054311 (39) | 0.889264 (36) |

In Table A6, we present the results obtained by the top-10 participants at the test phase. Similarly as in the previous phase, results are still very good with even lower bias scores, at least for the top participants, suggesting that participants were able to further improve their methods after the end of development phase. Another important aspect that can be seen is that, compared to the development

---

[6] The full leaderboards for both phases are shown in the supplementary material.

**Table 3.** Top-10 solutions on the test phase (and Baseline results). Top-3 winning solutions highlighted in bold. The number inside the parenthesis indicate the global rank position for that particular variable, used to compute the average ranking.

| Participant | Average Ranking | Entries | Bias (+ pairs) | Bias (- pairs) | Accuracy |
|---|---|---|---|---|---|
| **paranoidai** | 1.333333 (1) | 39 | 0.000059 (2) | 0.000012 (1) | 0.999966 (1) |
| **ustc-nelslip** | 3.666667 (2) | 12 | 0.000175 (4) | 0.000172 (2) | 0.999569 (5) |
| **CdtQin** | 4.000000 (3) | 25 | 0.000036 (1) | 0.000405 (9) | 0.999827 (2) |
| debias | 4.666667 (4) | 5 | 0.000036 (1) | 0.000460 (10) | 0.999825 (3) |
| zhaixingzi | 5.000000 (5) | 14 | 0.000116 (3) | 0.000237 (8) | 0.999698 (4) |
| bestone | 5.333333 (6) | 11 | 0.000175 (4) | 0.000197 (5) | 0.999565 (7) |
| haoxl | 5.333333 (6) | 31 | 0.000178 (6) | 0.000195 (4) | 0.999568 (6) |
| Early | 5.333333 (6) | 4 | 0.000175 (4) | 0.000190 (3) | 0.999547 (9) |
| lemoner20 | 7.000000 (7) | 9 | 0.000176 (5) | 0.000201 (6) | 0.999507 (10) |
| ai | 7.333333 (8) | 14 | 0.000180 (7) | 0.000217 (7) | 0.999560 (8) |
| *Baseline* | 34.666667 (28) | 3 | 0.059694 (33) | 0.058601 (36) | 0.859175 (35) |

phase, participants made an overall smaller number of submission, which can be explained due to two main reasons: 1) they had around 1 week to make submissions to the test phase (to avoid cheating related issues, also verified at the code verification stage, as they would have access to the test data, i.e., without labels); 2) we fixed the maximum number of submissions per day to 5 to avoid participants to improve the results on the test set by try and error.

## 5.2   Top Winning Approaches

This section briefly presents the top-winning approaches (shown in Table A6), specially those that agreed to share with the organizers the code (verified at the code verification stage) and fact sheets (containing detailed information about their methods), according to the rules of the challenge. Table 4 shows some general information about the top-3 winning approaches. The workflow diagrams of top-3 winning solutions are shown in the supplementary material.

**Table 4.** General information about the top-3 winning approaches.

| Features / Team | 1st: **paranoidai** | 2nd: **ustc-nelslip** | 3rd: **CdtQin** |
|---|---|---|---|
| Pre-trained models | - | √ | √ |
| External data | √ | √ | √ |
| Regularization strategies | - | √ | √ |
| Handcrafted features | - | - | - |
| Face detection, alignment or segmentation strategy | √ | √ | √ |
| Ensemble models | √ | √ | - |
| Different models for different protected groups | - | - | - |
| Explicitly classify the legitimate attributes | - | - | - |
| Explicitly classify other attributes (e.g., image quality) | - | - | - |
| Pre-processing bias mitigation (e.g. rebalancing training data) | - | √ | √ |
| In-processing bias mitigation (e.g. bias aware loss function) | √ | - | √ |
| Post-processing bias mitigation technique | √ | - | √ |

**1st place:** *paranoidai* [**70**][7] team proposed an asymmetric-arc-loss training and multi-step fine-tuning. Their motivation was based on observation that even two different people have typically some similarity, and trying to minimise such similarity may make the model pay useless attention to easy negative samples. To address this problem, they alter the convergence target such that easy negative samples contribute less to the final gradient. They first train a general model (ResNet101 as backbone) and perform its multi-step fine-tuning. To improve the performance they also employ several tricks such as re-ranking, boundary cut and hard-sample model fusion. According to them, the hard-sample model fusion significantly helped to mitigate bias. For this, they assume that after getting a final model, there must be some data on the training set that the model cannot predict correctly. These are obvious hard samples. To address this problem, they propose a model fusion strategy, where a fine-tuned model is built for false-positive results, in addition to another model which performs better for those hard samples but worse in general cases. At the fusion step, they only take the result with extremely high confidence from the hard-sample model.

**2nd place:** *ustc-nelslip* [**67**][8] team addressed the problem focusing on data balancing and ensemble models. First, they tested different face detection algorithms to find an effective face cropped method [35]. Then, a data re-sampling method is used to balance the data distribution by under-sampling the majority class (based on gender and skin colour), combined with the use of external data. Next, different training data enhancement methods are used to increase the diversity of samples by means of image quality and light conditions, for instance, with the goal to improve performance. Finally, the prediction results of eight different models having different backbones (ResNet50 and ResNet152) and head loss (e.g., Arcface [17] and Cosface [62]) are linearly combined at test stage.

**3rd place:** *CdtQin*[9] team presented a multi-branch training approach, using a modified ResNet-101 as backbone, with similarity distribution constraints. The similarity distributions for these branches are estimated and constrained, with the goal of forcing the same kind of distribution among different groups to be closer and the distance between positive and negative distributions to be larger. To this end, hard positive pairs are defined offline, while top-k hard negative pairs are selected online for each branch. The cosine similarity of these pairs is computed, and the estimated distribution is obtained as in [28]. For the drawn distributions, three constrains are considered, specifically *kl_loss*, *order_loss* and *entropy_loss*. The first measures the KL Divergence of two different groups (e.g., females with dark *vs.* light skin colour). The *order_loss* measures the expected difference with respect to two distributions. Intuitively, it is desired a large margin between positive and negative distributions. So, this loss is applied on the positive and negative similarity distributions for each branch. Finally, *entropy_loss*

---

[7] https://github.com/paranoidai/Fairface-Recognition-Solution

[8] https://github.com/HaoSir/ECCV-2020-Fair-Face-Recognition-challenge_2nd_place_solution-ustc-nelslip-

[9] https://github.com/CdtQin/FairFace

measures the negative entropy of a single distribution, designed to allows the similarity distribution near the threshold to have lower variance, promoting a better separation. The final loss is defined by a linear combination of these losses in addition to the ArcFace Loss [17].

### 5.3   Bias Analysis

In this section we analyze biases in the results of top-10 teams and discuss their possible causes. To conduct the analysis, we removed from the test set two error images found after the test stage was closed (one non-face, one wrong identity), which reduced the number of positive matches by 56 but did not affect the number of combinations of legitimate attributes. The changes to the calculated values of bias and accuracy were therefore very small and did not affect the findings nor changed the ranking of the top-3 teams. Detailed analyses are provided in the supplementary material. Main findings are summarized next.

**Breakdown of Average Discrimination:** A discrimination $d$, as defined by Eq. 3, quantifies the difference in accuracy between a given protected group and the best achieved one. High average discrimination of certain protected group therefore indicates that the accuracy of the algorithm is lower than for other protected groups. The character of bias we found in the algorithms of the top teams was not that they would have higher accuracy in all circumstances for certain groups and lower for others, but instead that they consider people from certain groups more similar to each other that individuals from other groups. Specifically, even though the differences were small, the algorithms consistently had difficulties distinguishing females with dark skin colour. This resulted in the lowest values of discrimination in the positive samples and the highest in the negative ones. Considering the averages over the top-10 teams, in positive samples the group with the highest discrimination were males with dark skin colour: $d = 4.748\text{e-}04$ (males with white skin colour were very close with $d = 4.690\text{e-}04$) and females with dark skin colour were the least discriminated: $d = 2.349\text{e-}04$. Conversely, in the negative samples females with dark skin colour were the most discriminated: $d = 1.783\text{e-}04$ and males with light skin colour least with $d = 0.475\text{e-}04$. Note however, that there were some exceptions from this trend. For example, for team *paranoidai* the least discriminated group in positive samples were not females with dark skin colour, but males with dark skin colour.

In addition to the absolute values of discrimination we also calculated for each protected group the frequency how often it was the most discriminated one (over all combinations of legitimate attributes). Even if a group is the most discriminated in 100% of the cases, the actual differences from the other groups might still be negligible. Nevertheless, it is convenient for showing trends as it allows to filter out outliers For the top-10 teams, in positive samples males with light skin colour were the most often discriminated group (42.2% cases) whereas females with dark skin colour the least often (11.2%). This was almost perfectly reversed in negative samples: females with dark skin colour were the

most frequently the group with the highest discrimination value (45.5%) whereas males with light skin colour were the least often group (12.6%). The exception was *paranoidai*, with the lowest frequency for females with dark skin colour in both positive and negative samples.

**Impact of Legitimate Attributes on Average Discrimination:** To analyze the effect of legitimate attributes we split their combinations into as many subgroups as there are possible values of the chosen attribute. For example, for glasses there are three subsets, the first one contains all samples where none of the images contain glasses, second group consists of samples where both images contain glasses and the third group are the remaining ones. We found that for some teams, wearing glasses makes individuals in both positive and negative samples look more similar in the sense that the differences in accuracy in positive samples tend to be the smallest if both images contain glasses and in negative samples the largest (note that in positive samples, teams *ustc-nelslip*, *bestone*, *haoxl*, *ai* are exceptions from this observation but in negative samples it holds for all top-10 teams). This is to a large extent an expected result: glasses cover part of the face which is one of the most important for recognition and therefore make people look more similar to each other.

Age was another attribute that clearly influenced the magnitude of bias: all top-10 teams exhibited higher values of discrimination in positive samples where both individuals were younger than 35 years and for the majority of the teams this was reversed in the negative samples, where the largest differences were obtained for the oldest subset (both individuals older than 65 years; exceptions are teams *paranoidai*, *CdtQin* and *debias*). This corresponds to findings of [50] and those from FRVT test [23] (which however emphasizes frequent exceptions). By analyzing the results further we found that in the competition dataset young individuals are less likely to wear glasses than the older ones. When considering only combinations of legitimate attributes where both individuals are younger than 35 years, only in 16% of them both individuals wear glasses but this ratio increases to 27.3% and 53.23% for the middle age and old subsets. Given the findings we made for the glasses attribute it is conceivable, that these two attributes act as magnifiers for each other.

Furthermore, we analyzed the effect of the remaining three legitimate attributes, i.e., head pose, image source and bounding box size. We did not find any clear trends shared by majority of the top-10 teams.

**Hardest Samples:** Hardest samples for top-3 teams are shown in Fig. 4. Even though the samples are different for different teams, they share common characteristics. The hardest positive samples are often composed from one "normal" image and one with extreme head pose or appearance variation, which makes them look differently. In the hardest negative samples on the other hand both images have often extreme head pose or glasses, which obscure parts of the faces important for the recognition and makes them look similar to each other.
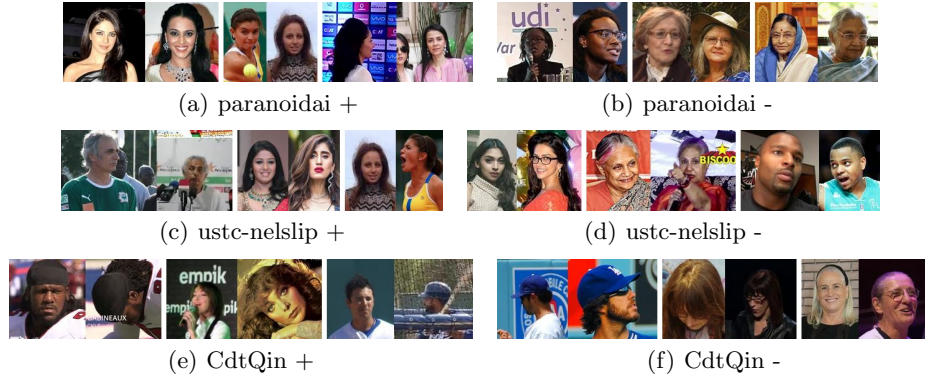
(a) paranoidai +                     (b) paranoidai -

(c) ustc-nelslip +                   (d) ustc-nelslip -

(e) CdtQin +                         (f) CdtQin -

**Fig. 4.** Most difficult samples for the top-3 teams: positive samples with lowest score (+), negative samples with highest score (-).

## Acknowledgments

## 6    Conclusions

This work presented the design and results of the FairFace Recognition Challenge at ECCV'2020. The submissions were evaluated on a reannotated version of IJB-C [37] database enriched by newly collected 12,549 public domain images. The participants were ranked using a novel evaluation protocol where both accuracy and bias scores were considered. The challenge attracted 151 participants. Top winning solutions obtained high performance in terms of accuracy ($\geq 0.999$ AUC-ROC) and bias scores. The post challenge analysis showed that top winning solutions applied a combination of different strategies to mitigate bias, such as face pre-processing, homogenization of data distributions, the use of bias aware loss functions and ensemble models, among others, suggesting there is not a general approach that works better for all the cases. Despite the high accuracy none of the methods was free of bias. By analysing the results of top-10 teams we found that their algorithms tend to have higher false positive rates for females with dark skin tone and for samples where both individuals wear glasses. In contrast there were higher false negative rates for males with light skin tone and for samples where both individuals are younger than 35 years. We also found that in the dataset individuals younger than 35 years wear glasses less often than older individuals, resulting in a combination of effects of these attributes.

# References

1. Facial recognition tech under spotlight after boston bombings. Biometric Technology Today **2013**(5), 1 (2013)
2. Albiero, V., Bowyer, K.W., Vangara, K., King, M.C.: Does face recognition accuracy get better with age? Deep face matchers say no. Winter Conference on Applications of Computer Vision (WACV) pp. 250–258 (2020)
3. Albiero, V., S., K.K., Vangara, K., Zhang, K., King, M.C., Bowyer, K.W.: Analysis of gender inequality in face recognition accuracy. CoRR **abs/2002.00065** (2020)
4. Alvi, M.S., Zisserman, A., Nellåker, C.: Turning a blind eye: Explicit removal of biases and variation from deep neural network embeddings. CoRR **abs/1809.02169** (2018)
5. American Civil Liberties Union, by Jacob Snow: Amazon's Face Recognition Falsely Matched 28 Members of Congress With Mugshots. [online] Available at: https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28 [Accessed 5 Aug. 2020] (July 2018)
6. Anne Hendricks, L., Burns, K., Saenko, K., Darrell, T., Rohrbach, A.: Women also snowboard: Overcoming bias in captioning models. In: European Conference on Computer Vision (ECCV). pp. 793–811 (2018)
7. Bellamy, R.K.E., Dey, K., Hind, M., Hoffman, S.C., Houde, S., Kannan, K., Lohia, P., Martino, J., Mehta, S., Mojsilovic, A., Nagar, S., Ramamurthy, K.N., Richards, J.T., Saha, D., Sattigeri, P., Singh, M., Varshney, K.R., Zhang, Y.: AI fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. CoRR **abs/1810.01943** (2018)
8. Berk, R., Heidari, H., Jabbari, S., Kearns, M., Roth, A.: Fairness in criminal justice risk assessments: The state of the art. Sociological Methods & Research (2018)
9. Bino, S., Bernerd, F.: Variations in skin colour and the biological consequences of ultraviolet radiation exposure. British Journal of Dermatology **169**(s3), 33–40 (2013)
10. Bird, S., Hutchinson, B., Kenthapadi, K., Kıcıman, E., Mitchell, M.: Fairness-aware machine learning: Practical challenges and lessons learned. In: Companion Proceedings of The 2019 World Wide Web Conference. p. 1297–1298 (2019)
11. Buolamwini, J., Gebru, T.: Gender shades: Intersectional accuracy disparities in commercial gender classification. In: Proceedings of the 1st Conference on Fairness, Accountability and Transparency. Proceedings of Machine Learning Research, vol. 81, pp. 77–91. PMLR (2018)
12. Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A.: VGGFace2: A dataset for recognising faces across pose and age. In: International Conference on Automatic Face Gesture Recognition (FG). pp. 67–74 (2018)
13. Cavazos, J.G., Phillips, P.J., Castillo, C.D., O'Toole, A.J.: Accuracy comparison across face recognition algorithms: Where are we on measuring race bias? CoRR **abs/1912.07398** (2019)
14. Chopra, S., Hadsell, R., LeCun, Y.: Learning a similarity metric discriminatively, with application to face verification. In: Conference on Computer Vision and Pattern Recognition (CVPR). vol. 1, pp. 539–546 (2005)
15. Chouldechova, A.: Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. Big Data **5**(2), 153–163 (2017)
16. Davies, B., Innes, M., Dawson, A.: An Evaluation of South Wales Police's Use of Automated Facial Recognition. [online] Available at: https://static1.squarespace.com/static/51b06364e4b02de2f57fd72e/t/5bfd4fbc21c67c2cdd692fa8/

1543327693640/AFR+Report+%5BDigital%5D.pdf [Accessed 5 Aug. 2020] (Sep 2018)

17. Deng, J., Guo, J., Xue, N., Zafeiriou, S.: ArcFace: Additive angular margin loss for deep face recognition. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

18. Dieterich, W., Mendoza, C., Brennan, T.: Compas risk scales : Demonstrating accuracy equity and predictive parity performance of the compas risk scales in broward county. [online] Available at: https://go.volarisgroup.com/ rs/430-MBX-989/images/ProPublica_Commentary_Final_070616.pdf [Accessed 5 Aug. 2020] (July 2016)

19. Drozdowski, P., Rathgeb, C., Dantcheva, A., Damer, N., Busch, C.: Demographic bias in biometrics: A survey on an emerging challenge. CoRR **abs/2003.02488** (2020)

20. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.S.: Fairness through awareness. CoRR **abs/1104.3913** (2011)

21. Escalante, H.J., Kaya, H., Salah, A., Escalera, S., Güçlütürk, Y., Güçlü, U., Baró, X., Guyon, I., Jacques Junior, J.C.S., Madadi, M., Ayache, S., Viegas, E., Gurpinar, F., Wicaksana, A.S., Liem, C., Van Gerven, M.A.J., Van Lier, R.: Modeling, recognizing, and explaining apparent personality from videos. IEEE Transactions on Affective Computing (2020)

22. Friedler, S.A., Scheidegger, C., Venkatasubramanian, S., Choudhary, S., Hamilton, E.P., Roth, D.: A comparative study of fairness-enhancing interventions in machine learning. In: Proceedings of the Conference on Fairness, Accountability, and Transparency. pp. 329–338. ACM (2019)

23. Grother, P., Ngan, M., Hanaoka, K.: Face Recognition Vendor Test (FRVT) Part 3: Demographic Effects. Tech. rep., National Institute of Standards and Technology (NIST) Interagency/Internal Report (NISTIR) - 8280 (2019)

24. Guo, G., Zhang, N.: A survey on deep learning based face recognition. Computer Vision and Image Understanding **189**, 102805 (2019)

25. Guo, Y., Zhang, L., Hu, Y., He, X., Gao, J.: MS-Celeb-1M: A dataset and benchmark for large-scale face recognition. CoRR **abs/1607.08221** (2016)

26. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778 (2016)

27. Huang, C., Li, Y., Loy, C.C., Tang, X.: Deep imbalanced learning for face recognition and attribute prediction. CoRR **abs/1806.00194** (2018)

28. Huang, Y., Shen, P., Tai, Y., Li, S., Liu, X., Li, J., Huang, F., Ji, R.: Improving face recognition from hard samples via distribution distillation loss. CoRR **abs/2002.03662** (2020)

29. Jayaraman, U., Gupta, P., Gupta, S., Arora, G., Tiwari, K.: Recent development in face recognition. Neurocomputing **408**, 231 – 245 (2020)

30. Kemelmacher-Shlizerman, I., Seitz, S.M., Miller, D., Brossard, E.: The megaface benchmark: 1 million faces for recognition at scale. In: Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4873–4882 (2016)

31. Kleinberg, J., Mullainathan, S., Raghavan, M.: Inherent trade-offs in the fair determination of risk scores. CoRR **abs/1609.05807** (2016)

32. Kortylewski, A., Egger, B., Schneider, A., Gerig, T., Morel-Forster, A., Vetter, T.: Analyzing and reducing the damage of dataset bias to face recognition with synthetic data. In: Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 1–8 (2019)

33. Kusner, M.J., Loftus, J., Russell, C., Silva, R.: Counterfactual fairness. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems 30, pp. 4066–4076. Curran Associates, Inc. (2017)
34. Learned-Miller, E., Huang, G.B., RoyChowdhury, A., Li, H., Hua, G.: Labeled Faces in the Wild: A Survey, pp. 189–248. Springer Publishing Company, Incorporated, 1st edn. (2016), Advances in Face Detection and Facial Image Analysis
35. Li, J., Wang, Y., Wang, C., Tai, Y., Qian, J., Yang, J., Wang, C., Li, J., Huang, F.: DSFD: Dual shot face detector. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
36. Lo Piano, S.: Ethical principles in machine learning and artificial intelligence: cases from the field and possible ways forward. Humanities and Social Sciences Communications **7**(9), 1–7 (2020)
37. Maze, B., Adams, J., Duncan, J.A., Kalka, N., Miller, T., Otto, C., Jain, A.K., Niggel, W.T., Anderson, J., Cheney, J., Grother, P.: IARPA Janus Benchmark - C: Face dataset and protocol. In: International Conference on Biometrics (ICB). pp. 158–165 (2018)
38. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. CoRR **abs/1908.09635** (2019)
39. Merler, M., Ratha, N.K., Feris, R.S., Smith, J.R.: Diversity in faces. CoRR **abs/1901.10436** (2019)
40. Morales, A., Fiérrez, J., Vera-Rodríguez, R.: Sensitivenets: Learning agnostic representations with application to face recognition. CoRR **abs/1902.00334** (2019)
41. Nabi, R., Shpitser, I.: Fair inference on outcomes. CoRR **abs/1705.10378** (2017)
42. Pearl, J.: Causal inference in statistics: An overview. Statistics Surveys **3**, 96–146 (2009)
43. Pessach, D., Shmueli, E.: Algorithmic fairness. CoRR **abs/2001.09784** (2020)
44. Pierce, J., Wong, R.Y., Merrill, N.: Sensor illumination: Exploring design qualities and ethical implications of smart cameras and image/video analytics. In: Conference on Human Factors in Computing Systems. p. 1–19 (2020)
45. ProPublica, by Julia Angwin, Larson, J., Mattu, S., Kirchner, L.: Machine bias: There's software used across the country to predict future criminals and it's biased against blacks. [online] Available at: https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing [Accessed 5 Aug. 2020] (May 2016)
46. Raji, I.D., Gebru, T., Mitchell, M., Buolamwini, J., Lee, J., Denton, E.: Saving face: Investigating the ethical concerns of facial recognition auditing. In: Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society. p. 145–151 (2020)
47. Robinson, J.P., Livitz, G., Henon, Y., Qin, C., Fu, Y., Timoner, S.: Face recognition: Too bias, or not too bias? In: Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 1–10 (2020)
48. Rothe, R., Timofte, R., Gool, L.V.: DEX: Deep expectation of apparent age from a single image. In: International Conference on Computer Vision Workshops (ICCVW). pp. 252–257 (2015)
49. Rothe, R., Timofte, R., Gool, L.V.: Deep expectation of real and apparent age from a single image without facial landmarks. International Journal of Computer Vision **126**(2-4), 144–157 (2018)
50. Srinivas, N., Ricanek, K., Michalski, D., Bolme, D.S., King, M.: Face recognition algorithm bias: Performance differences on images of children and adults. In: Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 2269–2277 (2019)

51. Terhörst, P., Kolf, J.N., Damer, N., Kirchbuchner, F., Kuijper, A.: Post-comparison mitigation of demographic bias in face recognition using fair score normalization. CoRR **abs/2002.03592** (2020)
52. The Guardian, by Wang Xueying: China testing facial-recognition surveillance system in Xinjiang - report. [online] Available at: https://www.theguardian.com/world/2018/jan/18/china-testing-facial-recognition-surveillance-system-in-xinjiang-report [Accessed 5 Aug. 2020] (Jan 2018)
53. The New York Times, by Jennifer Valentino-DeVries: How the Police Use Facial Recognition, and Where It Falls Short. [online] Available at: https://www.nytimes.com/2020/01/12/technology/facial-recognition-police.html [Accessed 5 Aug. 2020] (January 2020)
54. The New York Times, by Paul Mozur: Inside China's Dystopian Dreams: A.I., Shame and Lots of Cameras. [online] Available at: https://www.nytimes.com/2018/07/08/business/china-surveillance-technology.html [Accessed 5 Aug. 2020] (July 2018)
55. The New York Times, by Paul Mozur: One Month, 500,000 Face Scans: How China Is Using A.I. to Profile a Minority. [online] Available at: https://www.nytimes.com/2019/04/14/technology/china-surveillance-artificial-intelligence-racial-profiling.html [Accessed 5 Aug. 2020] (April 2019)
56. The Washington Post, by Jay Greene: Microsoft won't sell police its facial-recognition technology, following similar moves by Amazon and IBM. [online] Available at: https://www.washingtonpost.com/technology/2020/06/11/microsoft-facial-recognition [Accessed 5 Aug. 2020] (June 2020)
57. THINKPolicy Blog, by Arvind Krishna: IBM CEO's Letter to Congress on Racial Justice Reform. [online] Available at: https://www.ibm.com/blogs/policy/facial-recognition-sunset-racial-justice-reforms [Accessed 5 Aug. 2020] (June 2020)
58. Torralba, A., Efros, A.A.: Unbiased look at dataset bias. In: Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1521–1528 (2011)
59. US Day One Blog: We are implementing a one-year moratorium on police use of rekognition. [online] Available at: https://blog.aboutamazon.com/policy/we-are-implementing-a-one-year-moratorium-on-police-use-of-rekognition [Accessed 5 Aug. 2020] (June 2020)
60. Verma, S., Rubin, J.: Fairness definitions explained. In: Proceedings of the International Workshop on Software Fairness. pp. 1–7 (2018)
61. Vowels, M.J., Camgoz, N.C., Bowden, R.: NestedVAE: Isolating common factors via weak supervision. In: Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9202–9212 (2020)
62. Wang, H., Wang, Y., Zhou, Z., Ji, X., Gong, D., Zhou, J., Li, Z., Liu, W.: CosFace: Large margin cosine loss for deep face recognition. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
63. Wang, M., Deng, W.: Mitigating bias in face recognition using skewness-aware reinforcement learning. In: Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9322–9331 (2020)
64. Wang, M., Deng, W., Hu, J., Tao, X., Huang, Y.: Racial faces in-the-wild: Reducing racial bias by information maximization adaptation network. CoRR **abs/1812.00194** (2018)
65. Wang, T., Zhao, J., Yatskar, M., Chang, K.W., Ordonez, V.: Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In: International Conference on Computer Vision (ICCV). pp. 5310–5319 (2019)

66. Wang, Z., Qinami, K., Karakozis, I.C., Genova, K., Nair, P., Hata, K., Russakovsky, O.: Towards fairness in visual recognition: Effective strategies for bias mitigation. In: Conference on Computer Vision and Pattern Recognition (CVPR). pp. 8919–8928 (2020)
67. Yu, J., Hao, X., Xie, H., Yu, Y.: Fair face recognition using data balancing, enhancement and fusion. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (ECCVW), *in press* (2020)
68. Yucer, S., Akcay, S., Al-Moubayed, N., Breckon, T.P.: Exploring racial bias within face recognition via per-subject adversarially-enabled data augmentation. In: Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) (2020)
69. Zhang, L., Wu, Y., Wu, X.: A causal framework for discovering and removing direct and indirect discrimination. CoRR **abs/1611.07509** (2016)
70. Zhou, S.: AsArcFace: Asymmetric additive angular margin loss for fairface recognition. In: Proceedings of the European Conference on Computer Vision Workshops (ECCVW), *in press* (2020)

## Appendix - *Supplementary material*

## A    Introduction

This is the supplementary material for FairFace Challenge at ECCV 2020: Analyzing Bias in Face Recognition, a summary paper of the 2020 ChaLearn Looking at People Fair Face Recognition and Analysis Challenge held at ECCV 2020. Sec. B describes detailed schedule of the challenge, Sec. C its general statistics, Sec. D shows final leaderboards of both development and test phases, Sec. E shows workflows of the top-3 methods, Sec. F contains source tables for the Bias Analysis section in the main paper and Sec. G summarizes the instructions given to the annotators.

## B    Schedule

The schedule of the competition was as follows:

- **Apr 4th, 2020**. Start of the Challenge (development phase) – Release of training (with ground truth) and validation data (without ground truth).
- **Jun 22th, 2020**. End of development phase / Start of test phase – Release of test data (without ground truth) and validation labels.
- **Jul 1st, 2020**. End of the Challenge – Deadline for submitting the final predictions over the test (evaluation) data.
- **Jul 7th, 2020**. Submission of code and fact sheets – Containing detailed instructions to reproduce the results obtained on the test set and fact sheets with detailed and technical information about the developed approach.
- **Jul 12th, 2020**. Release of final results (after code verification).

## C    General Statistics

Fig. A5 shows the number of submissions per day on each phase, where a higher activity can be observed close to the end of each phase, indicating that participants may be fine tuning their methods and making more submissions in order to improve their rank positions.

## D    Full Leaderboard: development and test phase

The complete leaderboard of the development and test phases are shown in Table A5 and Table A6, respectively, for participants showing accuracy higher than 80%.

## E    Workflow of Top-3 Winning Approaches

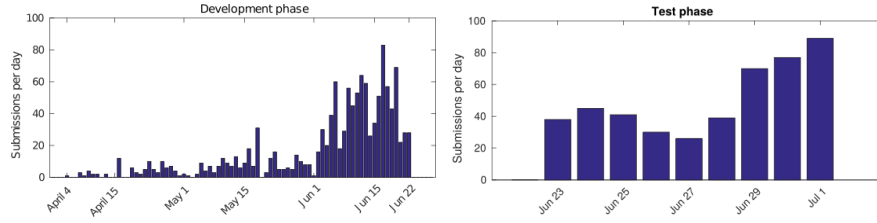The workflow diagram of top-3 winning solutions is shown in Fig. A6.

**Fig. A5.** Challenge evolution: number of submissions per day.

**Table A5.** Leaderboard of the development phase. The number inside the parenthesis indicate the global rank position for that particular variable, used to compute the average ranking.

| Parcitipant | Average Ranking | Entries | Bias (+ pairs) | Bias (- pairs) | Accuracy |
|---|---|---|---|---|---|
| ustc-nelslip | 2.333333 (1) | 30 | 0.000142 (1) | 0.002956 (3) | 0.999287 (3) |
| zheng.zhu | 3.666667 (2) | 133 | 0.000344 (3) | 0.003781 (7) | 0.999442 (1) |
| CdtQin | 3.666667 (2) | 72 | 0.000472 (5) | 0.002334 (1) | 0.998477 (5) |
| crisp | 4.666667 (3) | 14 | 0.000935 (8) | 0.003193 (4) | 0.999394 (2) |
| haoxl | 4.666667 (3) | 73 | 0.000348 (4) | 0.003678 (6) | 0.998699 (4) |
| cam_vision | 5.000000 (4) | 95 | 0.000731 (6) | 0.002488 (2) | 0.995621 (7) |
| Hyg | 6.000000 (5) | 33 | 0.000814 (7) | 0.003305 (5) | 0.998402 (6) |
| senlin11 | 9.333333 (6) | 50 | 0.000165 (2) | 0.010091 (16) | 0.992093 (10) |
| hanamichi | 10.666667 (7) | 91 | 0.001631 (9) | 0.006760 (10) | 0.987382 (13) |
| paranoidai | 12.000000 (8) | 156 | 0.003779 (12) | 0.007745 (13) | 0.988359 (11) |
| six_god | 13.000000 (9) | 2 | 0.006400 (23) | 0.004670 (8) | 0.993343 (8) |
| vuvko | 13.333333 (10) | 10 | 0.005018 (17) | 0.006880 (11) | 0.988125 (12) |
| camel | 14.333333 (11) | 66 | 0.007078 (25) | 0.005986 (9) | 0.993202 (9) |
| debias | 15.333333 (12) | 23 | 0.002808 (11) | 0.010383 (17) | 0.977708 (18) |
| UAM_Ignacio | 15.666667 (13) | 59 | 0.005009 (16) | 0.010054 (15) | 0.981019 (16) |
| zhaixingzi | 15.666667 (13) | 61 | 0.005322 (19) | 0.010022 (14) | 0.984689 (14) |
| clessvna | 17.000000 (14) | 6 | 0.006617 (24) | 0.007572 (12) | 0.981362 (15) |
| dddddddqiu | 18.666667 (15) | 1 | 0.002675 (10) | 0.015141 (21) | 0.967389 (25) |
| jjjjjjjm | 19.000000 (16) | 1 | 0.004937 (15) | 0.013108 (19) | 0.972278 (23) |
| clearlove10 | 19.333333 (17) | 1 | 0.005039 (18) | 0.012280 (18) | 0.972329 (22) |
| ai | 19.666667 (18) | 2 | 0.004123 (13) | 0.019420 (25) | 0.974442 (21) |
| hanhao1415 | 20.000000 (19) | 13 | 0.005686 (21) | 0.014742 (20) | 0.977388 (19) |
| zhangkun | 21.666667 (20) | 7 | 0.004896 (14) | 0.020322 (27) | 0.968208 (24) |
| YSTBER | 23.000000 (21) | 1 | 0.008320 (27) | 0.015763 (22) | 0.977343 (20) |
| season | 24.000000 (22) | 4 | 0.011135 (31) | 0.017689 (24) | 0.978085 (17) |
| TCxu | 25.333333 (23) | 33 | 0.005972 (22) | 0.021329 (28) | 0.964486 (26) |
| Finn_zhang | 28.333333 (24) | 15 | 0.005468 (20) | 0.044931 (36) | 0.947747 (29) |
| okpeng | 29.000000 (25) | 42 | 0.010581 (30) | 0.025169 (30) | 0.949839 (27) |
| wg1234567p | 30.666667 (26) | 14 | 0.007765 (26) | 0.042682 (35) | 0.939736 (31) |
| Serendi | 31.000000 (27) | 5 | 0.021709 (34) | 0.021855 (29) | 0.946445 (30) |
| baoqianyue | 31.000000 (27) | 44 | 0.009132 (28) | 0.031606 (32) | 0.938430 (33) |
| suhk | 31.333333 (28) | 10 | 0.014761 (33) | 0.016355 (23) | 0.840568 (38) |
| burning | 32.333333 (29) | 52 | 0.024901 (37) | 0.020274 (26) | 0.915005 (34) |
| quentinyq | 32.333333 (29) | 1 | 0.022878 (35) | 0.037616 (34) | 0.949130 (28) |
| jieson_zheng | 33.666667 (30) | 1 | 0.014731 (32) | 0.045830 (37) | 0.939732 (32) |
| yuchun_wang | 34.666667 (31) | 14 | 0.010568 (29) | 0.052194 (38) | 0.868556 (37) |
| fireant | 34.666667 (31) | 1 | 0.023675 (36) | 0.034809 (33) | 0.903129 (35) |
| mengtzu.chiu | 36.000000 (32) | 9 | 0.050556 (38) | 0.025790 (31) | 0.837854 (39) |
| *Baseline* | 38.333333 (33) | 1 | 0.057620 (40) | 0.054311 (39) | 0.889264 (36) |
| VisTeam | 39.666667 (34) | 59 | 0.054725 (39) | 0.061032 (40) | 0.820067 (40) |

(a) *paranoidai*
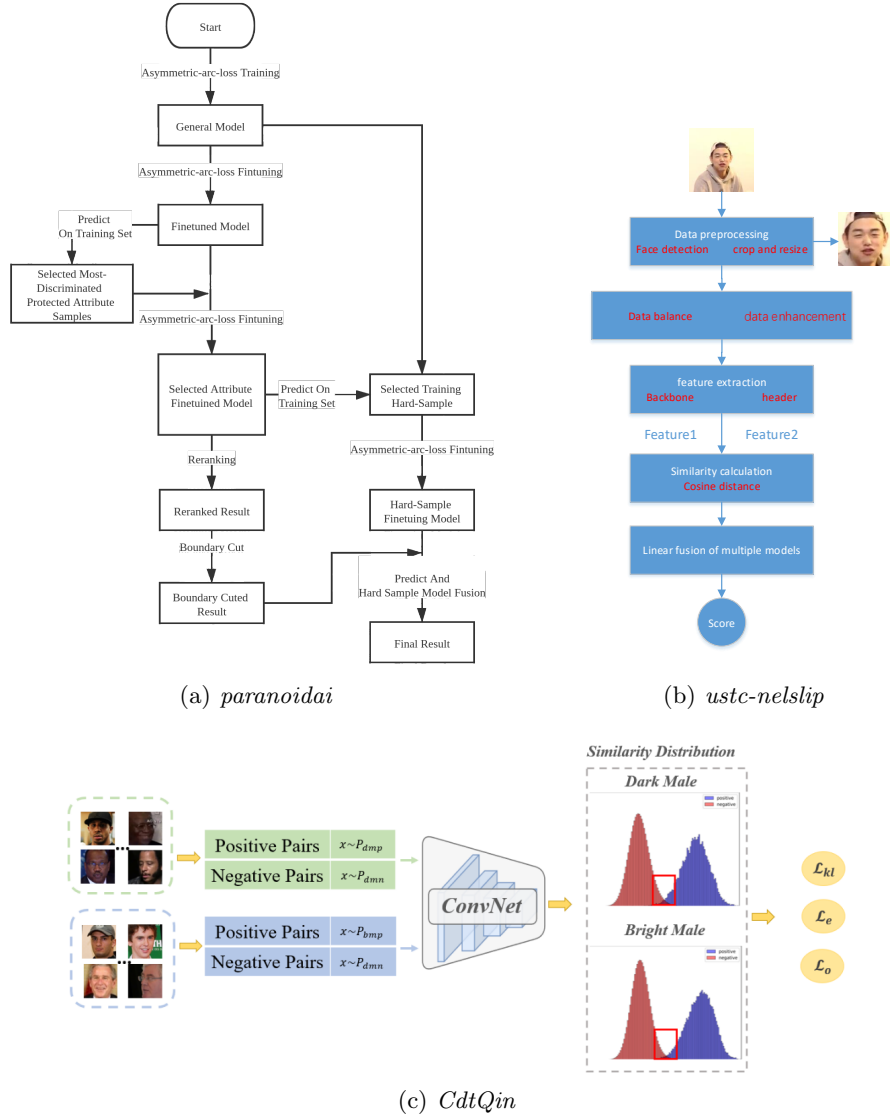


(b) *ustc-nelslip*



(c) *CdtQin*

**Fig. A6.** Workflow diagram of top-3 winning solutions.

**Table A6.** Leaderboard of the test phase. Top-3 winning solutions highlighted in bold. The number inside the parenthesis indicate the global rank position for that particular variable, used to compute the average ranking.

| Participant | Average Ranking | Entries | Bias (+ pairs) | Bias (- pairs) | Accuracy |
|---|---|---|---|---|---|
| **paranoidai** | 1.333333 (1) | 39 | 0.000059 (2) | 0.000012 (1) | 0.999966 (1) |
| **ustc-nelslip** | 3.666667 (2) | 12 | 0.000175 (4) | 0.000172 (2) | 0.999569 (5) |
| **CdtQin** | 4.000000 (3) | 25 | 0.000036 (1) | 0.000405 (9) | 0.999827 (2) |
| debias | 4.666667 (4) | 5 | 0.000036 (1) | 0.000460 (10) | 0.999825 (3) |
| zhaixingzi | 5.000000 (5) | 14 | 0.000116 (3) | 0.000237 (8) | 0.999698 (4) |
| bestone | 5.333333 (6) | 11 | 0.000175 (4) | 0.000197 (5) | 0.999565 (7) |
| haoxl | 5.333333 (6) | 31 | 0.000178 (6) | 0.000195 (4) | 0.999568 (6) |
| Early | 5.333333 (6) | 4 | 0.000175 (4) | 0.000190 (3) | 0.999547 (9) |
| lemoner20 | 7.000000 (7) | 9 | 0.000176 (5) | 0.000201 (6) | 0.999507 (10) |
| ai | 7.333333 (8) | 14 | 0.000180 (7) | 0.000217 (7) | 0.999560 (8) |
| six_god | 12.333333 (9) | 8 | 0.000540 (13) | 0.000984 (11) | 0.998785 (13) |
| mcga | 13.000000 (10) | 3 | 0.000341 (11) | 0.001228 (12) | 0.998265 (16) |
| lwx | 13.000000 (10) | 12 | 0.000327 (10) | 0.001444 (14) | 0.998545 (15) |
| doinb | 13.333333 (11) | 8 | 0.000580 (14) | 0.001599 (15) | 0.999297 (11) |
| clearlove10 | 13.333333 (11) | 4 | 0.000687 (15) | 0.001362 (13) | 0.999270 (12) |
| Hans | 13.666667 (12) | 24 | 0.000206 (8) | 0.002497 (16) | 0.998157 (17) |
| YSTBER | 14.333333 (13) | 10 | 0.000396 (12) | 0.003352 (17) | 0.998573 (14) |
| hanamichi | 15.000000 (14) | 11 | 0.000280 (9) | 0.005279 (18) | 0.996242 (18) |
| burning | 18.333333 (15) | 19 | 0.000969 (16) | 0.005815 (19) | 0.992119 (20) |
| zheng.zhu | 18.666667 (16) | 26 | 0.001206 (17) | 0.006573 (20) | 0.993509 (19) |
| hq2172 | 20.000000 (17) | 13 | 0.001503 (18) | 0.007151 (21) | 0.990733 (21) |
| vuvko | 22.000000 (18) | 10 | 0.003961 (21) | 0.007562 (22) | 0.983437 (23) |
| cam_vision | 22.000000 (18) | 21 | 0.002094 (19) | 0.008945 (25) | 0.989470 (22) |
| UAM_Ignacio | 23.000000 (19) | 21 | 0.003478 (20) | 0.008249 (23) | 0.974710 (26) |
| camel | 25.000000 (20) | 7 | 0.006143 (24) | 0.010392 (27) | 0.981795 (24) |
| DeepBlueAI | 25.333333 (21) | 5 | 0.008111 (25) | 0.009572 (26) | 0.977451 (25) |
| ztelily | 25.666667 (22) | 15 | 0.005236 (22) | 0.014847 (28) | 0.962481 (27) |
| baoqianyue | 27.666667 (23) | 3 | 0.005377 (23) | 0.021418 (31) | 0.951101 (29) |
| yuchun_wang | 28.666667 (24) | 11 | 0.011524 (28) | 0.008660 (24) | 0.881282 (34) |
| lijianshu | 28.666667 (24) | 3 | 0.008862 (26) | 0.021511 (32) | 0.962229 (28) |
| VisTeam | 31.000000 (25) | 15 | 0.019902 (31) | 0.016837 (29) | 0.917651 (33) |
| jieson_zheng | 31.000000 (25) | 4 | 0.011107 (27) | 0.033817 (35) | 0.941330 (31) |
| wg1234567p | 31.000000 (25) | 2 | 0.012173 (30) | 0.022290 (33) | 0.941947 (30) |
| Finn_zhang | 31.666667 (26) | 1 | 0.011554 (29) | 0.024265 (34) | 0.940516 (32) |
| mengtzu.chiu | 32.666667 (27) | 13 | 0.023490 (32) | 0.018914 (30) | 0.830624 (36) |
| *Baseline* | 34.666667 (28) | 3 | 0.059694 (33) | 0.058601 (36) | 0.859175 (35) |

# F   Bias Analysis

This section contains the source tables for the Bias Analysis section in the main paper.

## F.1   Breakdown of Average Discrimination

Table A7 shows average discriminations for every protected group. Table A8 shows frequencies, how often (over all combinations of legitimate attributes) was a given protected group the most discriminated one. Both tables contain results for the top-3 teams and averages for the top-10 teams. Values for positive samples are indicated by + after the team name and values for the negative samples by - after the team name.

**Table A7.** Average discrimination for top-3 teams and top-10 teams average. Each number is a mantissa $m$ in the scientific notation m.e-04.

| Participant | Male Light | Male Dark | Female Light | Female Dark |
|---|---|---|---|---|
| paranoidai + | 0.453 | 0.117 | 0.344 | 0.258 |
| paranoidai - | 0.166 | 0.138 | 0.194 | 0.762 |
| ustc-nelslip + | 5.491 | 5.747 | 5.208 | 3.217 |
| ustc-nelslip - | 0.555 | 0.961 | 1.308 | 2.306 |
| CdtQin + | 3.148 | 4.661 | 1.552 | 0.418 |
| CdtQin - | 0.373 | 0.552 | 0.413 | 0.740 |
| Top-10 avg + | 4.690 | 4.748 | 3.896 | 2.349 |
| Top-10 avg - | 0.475 | 0.775 | 0.982 | 1.783 |

**Table A8.** Frequency (over all combinations of legitimate attributes) of being the most discriminated protected group.

| Participant | Male Light | Male Dark | Female Light | Female Dark |
|---|---|---|---|---|
| paranoidai + | 0.401 | 0.260 | 0.249 | 0.090 |
| paranoidai - | 0.361 | 0.254 | 0.310 | 0.075 |
| ustc-nelslip + | 0.393 | 0.226 | 0.248 | 0.134 |
| ustc-nelslip - | 0.073 | 0.202 | 0.198 | 0.527 |
| CdtQin + | 0.391 | 0.350 | 0.170 | 0.090 |
| CdtQin - | 0.174 | 0.282 | 0.182 | 0.362 |
| Top-10 avg + | 0.422 | 0.246 | 0.220 | 0.112 |
| Top-10 avg - | 0.126 | 0.214 | 0.205 | 0.455 |

## F.2   Impact of Legitimate Attributes on Average Discrimination

Tables A9, A10, A11, A12, A13 and A14 demonstrate dependencies between average discrimination and attributes age, wearing glasses, head pose, bounding box size and image source respectively. All tables contain results for the top-3 teams and averages for the top-10 teams. Values for positive samples are indicated by + after the team name and values for the negative samples by - after the team name. For brevity, we denote the subset of the samples by using initial letter of the attribute followed by its label. For example, for glasses there are G0-G0 (no glasses in either of the images), G0-G1 (one image does not contain glassesm the other one does) and G1-G1 (both images contain glasses). For every subset we calculate average discrimination of each protected group but to save space we only report maximum and minimum values and denote the group in the parentheses: 1=Male Light, 2=Male Dark, 3=Female Light and 4=Female Dark. Each number is a mantissa $m$ in the scientific notation m.e-04.

**Table A9.** Effect of age (same age groups in the sample) on average discrimination.

| | max | | | min | | |
|---|---|---|---|---|---|---|
| Participant | A0-A0 | A1-A1 | A2-A2 | A0-A0 | A1-A1 | A2-A2 |
| paranoidai + | 1.693 (4) | 0.608 (3) | 0.162 (2) | 0.226 (2) | 0.009 (4) | 0.012 (4) |
| paranoidai + | 6.147 (4) | 0.162 (1) | 0.828 (3) | 0.234 (3) | 0.012 (4) | 0.179 (2) |
| ustc-nelslip + | 16.044 (1) | 3.562 (2) | 5.859 (2) | 1.426 (2) | 2.171 (4) | 0.5 (4) |
| ustc-nelslip - | 2.844 (4) | 1.615 (4) | 6.799 (4) | 1.15 (3) | 0.4 (1) | 0.724 (1) |
| CdtQin + | 14.547 (2) | 1.449 (2) | 2.492 (2) | 1.948 (4) | 0.078 (4) | 0.001 (4) |
| CdtQin - | 1.208 (4) | 0.54 (4) | 0.977 (4) | 0.396 (3) | 0.181 (1) | 0.39 (1) |
| Top-10 avg + | 13.389 (1) | 2.535 (2) | 3.931 (2) | 2.576 (3) | 1.615 (4) | 0.321 (4) |
| Top-10 avg - | 2.714 (4) | 1.233 (4) | 4.713 (4) | 0.891 (3) | 0.32 (1) | 0.591 (1) |

**Table A10.** Effect of age (different age groups in the sample) on average discrimination.

| | max | | | min | | |
|---|---|---|---|---|---|---|
| *Participant* | *A0-A1* | *A0-A2* | *A1-A2* | *A0-A1* | *A0-A2* | *A1-A2* |
| paranoidai + | 0.642 (1) | - | 0.469 (1) | 0.039 (4) | - | 0.011 (4) |
| paranoidai - | 0.166 (2) | 0.059 (1) | 0.116 (3) | 0.028 (4) | 0.006 (4) | 0.051 (2) |
| ustc-nelslip + | 16.684 (3) | - | 11.87 (2) | 1.217 (4) | - | 0.51 (4) |
| ustc-nelslip - | 1.181 (4) | 0.922 (4) | 2.591 (4) | 0.417 (1) | 0.583 (1) | 0.346 (1) |
| CdtQin + | 8.795 (2) | - | 0.795 (1) | 0.551 (4) | - | 0.003 (4) |
| CdtQin - | 0.703 (4) | 0.602 (2) | 0.681 (4) | 0.273 (3) | 0.409 (3) | 0.246 (1) |
| Top-10 avg + | 12.331 (3) | - | 7.484 (2) | 0.956 (4) | - | 0.476 (4) |
| Top-10 avg - | 0.939 (4) | 0.784 (4) | 1.945 (4) | 0.382 (1) | 0.534 (1) | 0.292 (1) |

**Table A11.** Effect of wearing glasses on average discrimination.

| | max | | | min | | |
|---|---|---|---|---|---|---|
| *Participant* | *G0-G0* | *G0-G1* | *G1-G1* | *G0-G0* | *G0-G1* | *G1-G1* |
| paranoidai + | 0.542 (1) | 0.838 (3) | 0.075 (1) | 0.146 (2) | 0.117 (2) | 0.015 (4) |
| paranoidai - | 0.176 (3) | 0.13 (1) | 2.924 (4) | 0.066 (4) | 0.038 (4) | 0.159 (2) |
| ustc-nelslip + | 9.098 (2) | 4.644 (3) | 7.01 (4) | 1.5 (4) | 1.407 (2) | 0.128 (3) |
| ustc-nelslip - | 2.242 (4) | 1.751 (4) | 3.478 (4) | 0.453 (1) | 0.439 (1) | 0.892 (1) |
| CdtQin + | 7.339 (2) | 8.956 (1) | 0.706 (1) | 0.62 (4) | 0.192 (4) | 0.067 (3) |
| CdtQin - | 0.5 (4) | 0.732 (4) | 1.004 (4) | 0.284 (1) | 0.345 (1) | 0.507 (3) |
| Top-10 avg + | 7.332 (2) | 5.617 (1) | 4.734 (4) | 1.252 (4) | 1.781 (4) | 0.121 (3) |
| Top-10 avg - | 1.579 (4) | 1.368 (4) | 2.82 (4) | 0.379 (1) | 0.388 (1) | 0.749 (1) |

**Table A12.** Effect of head pose on average discrimination.

| | max | | | min | | |
|---|---|---|---|---|---|---|
| *Participant* | *H0-H0* | *H0-H1* | *H1-H1* | *H0-H0* | *H0-H1* | *H1-H1* |
| paranoidai + | 0.794 (3) | 0.508 (1) | 0.284 (4) | 0.03 (2) | 0.116 (2) | 0.12 (3) |
| paranoidai - | 0.183 (3) | 0.362 (4) | 2.234 (4) | 0.05 (4) | 0.128 (2) | 0.166 (1) |
| ustc-nelslip + | 4.804 (3) | 9.473 (2) | 8.86 (4) | 0.493 (2) | 0.498 (4) | 4.055 (2) |
| ustc-nelslip - | 2.42 (4) | 2.279 (4) | 2.246 (4) | 0.605 (1) | 0.543 (1) | 0.527 (1) |
| CdtQin + | 2.898 (3) | 5.378 (2) | 8.158 (1) | 0.187 (2) | 0.214 (4) | 0.764 (4) |
| CdtQin - | 0.543 (2) | 0.735 (4) | 1.081 (4) | 0.234 (1) | 0.375 (1) | 0.502 (3) |
| Top-10 avg + | 3.601 (3) | 7.173 (2) | 7.808 (1) | 0.337 (2) | 0.464 (4) | 3.333 (3) |
| Top-10 avg - | 1.716 (4) | 1.729 (4) | 1.952 (4) | 0.474 (1) | 0.471 (1) | 0.485 (1) |

# G   Summary of Annotation Instructions

## G.1   Instructions for Annotators

- Gender: Use binary categories corresponding to biological sex: male and female.
- Skin colour: Compare the skin tone with provided templates for Fitzpatrick skin types (an example is shown in Fig. A7) and select the most similar one.

**Table A13.** Effect of bounding box size on average discrimination.

|  | max | | | min | | |
| --- | --- | --- | --- | --- | --- | --- |
| *Participant* | *B0-B0* | *B0-B1* | *B1-B1* | *B0-B0* | *B0-B1* | *B1-B1* |
| paranoidai + | 0.258 (3) | 0.456 (3) | 1.099 (1) | 0.006 (4) | 0.104 (2) | 0.183 (2) |
| paranoidai - | 0.172 (3) | 1.19 (4) | 0.623 (4) | 0.045 (4) | 0.122 (2) | 0.162 (2) |
| ustc-nelslip + | 6.035 (2) | 5.905 (3) | 9.611 (2) | 0.543 (4) | 2.473 (2) | 3.812 (4) |
| ustc-nelslip - | 2.127 (4) | 2.168 (4) | 2.741 (4) | 0.579 (1) | 0.562 (1) | 0.518 (1) |
| CdtQin + | 5.117 (2) | 5.352 (2) | 3.561 (2) | 0.001 (4) | 0.392 (4) | 0.248 (3) |
| CdtQin - | 0.729 (4) | 0.735 (4) | 0.77 (2) | 0.282 (3) | 0.349 (1) | 0.386 (1) |
| Top-10 avg + | 5.098 (2) | 4.321 (3) | 7.205 (2) | 0.371 (4) | 2.59 (2) | 3.049 (4) |
| Top-10 avg - | 1.618 (4) | 1.733 (4) | 2.037 (4) | 0.506 (1) | 0.474 (1) | 0.449 (1) |

**Table A14.** Effect of image source on average discrimination.

|  | max | | | min | | |
| --- | --- | --- | --- | --- | --- | --- |
| *Participant* | *S0-S0* | *S0-S1* | *S1-S1* | *S0-S0* | *S0-S1* | *S1-S1* |
| paranoidai + | 1.572 (1) | 0.265 (3) | 0.507 (3) | 0.151 (2) | 0.098 (2) | 0.012 (4) |
| paranoidai - | 0.32 (3) | 1.374 (4) | 0.156 (1) | 0.151 (1) | 0.127 (2) | 0.03 (4) |
| ustc-nelslip + | 11.568 (1) | 7.534 (2) | 4.274 (1) | 2.793 (4) | 3.662 (1) | 0.509 (2) |
| ustc-nelslip - | 2.833 (4) | 2.326 (4) | 1.753 (4) | 0.499 (1) | 0.512 (1) | 0.692 (1) |
| CdtQin + | 4.67 (3) | 5.643 (1) | 5.442 (2) | 0.373 (4) | 0.57 (4) | 0.09 (3) |
| CdtQin - | 0.908 (4) | 0.727 (4) | 0.604 (4) | 0.553 (1) | 0.342 (1) | 0.257 (1) |
| Top-10 avg + | 9.467 (1) | 5.96 (2) | 3.0 (1) | 2.08 (4) | 2.866 (4) | 1.254 (3) |
| Top-10 avg - | 2.058 (4) | 1.85 (4) | 1.385 (4) | 0.458 (1) | 0.446 (1) | 0.549 (1) |



**Fig. A7.** Example of a template for Fitzpatrick skin types given to the annotators. Available online at https://www.rejuvent.com/why-your-skin-type-is-important/ (accessed 8 Sep 2020).
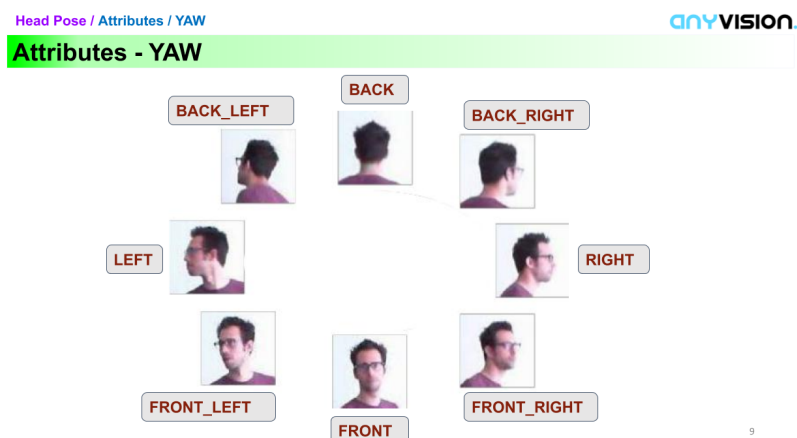
**Fig. A8.** Illustration of head yaw

- Age: Annotate perceived exact age.
- Head pose: Annotate head yaw as one of the 8 equally distributed directions as shown in Fig. A8.
- Glasses: Three labels, distinguish faces with no glasses, transparent glasses and sunglasses.
- Image source: Obtained automatically.
- Bounding box: Use original bounding box for IJB-C images, provide loose crop of the face for the newly collected ones.

### G.2    Label Aggregation & Post-processing

- Gender: Final labels synchronized for each identity to the most prevalent ones.
- Skin colour: two final categories: light corresponding to skin types I-III, dark corresponding to types IV-VI. Final labels synchronized for each identity to the most prevalent ones.
- Age group: Estimate of each annotator was adjusted by $age_{adj} = k_i \times age_{anno} + q_i$ (coefficients $k_i$ and $q_i$ were learned for each annotator by least squares from a subset of images for which the exact age was known). The final label for each image obtained by thresholding the mean of the adjusted estimates to three final categories: 0-34, 35-64 and 65+.
- Head pose: Two final categories: front, front left and front right are marked as 'frontal', other poses as 'other'. Final labels synchronized for each image to the most prevalent ones.
- Glasses: Labels for transparent and sunglasses were merged into a single category glasses. Final labels synchronized for each image to the most prevalent ones.

– Image source: None.
– Bounding box size: Two final categories: bounding boxes with both dimensions >224 px categorized as big, others as small.