
Robustness Disparities in Commercial Face Detection

Samuel Dooley
University of Maryland
sdooley1@cs.umd.edu

Tom Goldstein
University of Maryland
tomg@cs.umd.edu

John P. Dickerson
University of Maryland
john@cs.umd.edu

Abstract

1 Facial detection and analysis systems have been deployed by large companies and
2 critiqued by scholars and activists for the past decade. Critiques that focus on
3 system performance analyze disparity of the system’s output, i.e., how frequently is
4 a face detected for different Fitzpatrick skin types or perceived genders. However,
5 we focus on the robustness of these system outputs under noisy natural perturba-
6 tions. We present the first of its kind detailed benchmark of the robustness of two
7 such systems: Amazon Rekognition and Microsoft Azure. [We use both standard
8 and recently released academic facial datasets to quantitatively analyze trends in
9 robustness for each. Qualitatively across all the datasets and systems, we find that
10 photos of individuals who are *older, masculine presenting, of darker skin type, or
11 have dim lighting* are more susceptible to errors than their counterparts in other
12 identities.](#)

13 1 Introduction

14 Face detection systems identify the presence and location of faces in images and video. Automated
15 face detection is a core component of myriad systems—including *face recognition technologies*
16 (FRT), wherein a detected face is matched against a database of faces, typically for identification
17 or verification purposes. FRT-based systems are widely deployed [Hartzog, 2020, Derringer, 2019,
18 Weise and Singer, 2020]. Automated face recognition enables capabilities ranging from the relatively
19 morally neutral (e.g., searching for photos on a personal phone [Google, 2021]) to morally laden (e.g.,
20 widespread citizen surveillance [Hartzog, 2020], or target identification in warzones [Marson and
21 Forrest, 2021]). Legal and social norms regarding the usage of FRT are evolving [e.g., Grother et al.,
22 2019]. For example, in June 2021, the first county-wide ban on its use for policing [see, e.g., Garvie,
23 2016] went into effect in the US [Gutman, 2021]. Some use cases for FRT will be deemed socially
24 repugnant and thus be either legally or *de facto* banned from use; yet, it is likely that pervasive use of
25 facial analysis will remain—albeit with more guardrails than are found today [Singer, 2018].

26 One such guardrail that has spurred positive, though insufficient, improvements and widespread
27 attention is the use of benchmarks. For example, in late 2019, the US National Institute of Standards
28 and Technology (NIST) adapted its venerable Face Recognition Vendor Test (FRVT) to explicitly
29 include concerns for demographic effects [Grother et al., 2019], ensuring such concerns propagate
30 into industry systems. Yet, differential treatment by FRT of groups has been known for at least a
31 decade [e.g., Klare et al., 2012, El Khiyari and Wechsler, 2016], and more recent work spearheaded
32 by Buolamwini and Gebru [2018] uncovers unequal performance at the phenotypic subgroup level.
33 That latter work brought widespread public, and thus burgeoning regulatory, attention to bias in
34 FRT [e.g., Lohr, 2018, Kantayya, 2020].

35 One yet unexplored benchmark examines the bias present in a system’s robustness (e.g., to noise, or
36 to different lighting conditions), both in aggregate and with respect to different dimensions of the
37 population on which it will be used. Many detection and recognition systems are not built in house,
38 instead making use of commercial cloud-based “ML as a Service” (MLaaS) platforms offered by
39 tech giants such as Amazon and Microsoft. The implementation details of those systems are not

40 exposed to the end user—and even if they were, quantifying their failure modes would be difficult.
41 With this in mind, our **main contribution** is a wide *robustness benchmark* of two commercial-grade
42 face detection systems (accessed via Amazon’s Rekognition and Microsoft’s Azure face detection
43 APIs). For fifteen types of realistic noise, and five levels of severity per type of noise [Hendrycks and
44 Dietterich, 2019], we test both APIs against images in each of four well-known datasets. Across these
45 more than 5,000,000 noisy images, we analyze the impact of noise on face detection performance.
46 Perhaps unsurprisingly, we find that noise decreases overall performance, and that different types of
47 noise impact, in an “unfair” way, cross sections of the population of images (e.g., based on Fitzgerald
48 skin type, age, self-identified gender, and intersections of those dimensions). Our method is extensible
49 and can be used to quantify the robustness of other detection and FRT systems, and adds to the
50 burgeoning literature supporting the necessity of explicitly considering fairness in ML systems with
51 morally-laden downstream uses.

52 2 Related Work

53 We briefly overview additional related work in the two core areas addressed by our benchmark:
54 robustness to noise and demographic disparity in facial detection and recognition. That latter point
55 overlaps heavily with the fairness in machine learning literature; for additional coverage of that
56 broader ecosystem and discussion around fairness in machine learning writ large, we direct the reader
57 to survey works due to Chouldechova and Roth [2018] and Barocas et al. [2019].

58 **Demographic effects in facial detection and recognition.** The existence of differential perfor-
59 mance of facial detection and recognition on groups and subgroups of populations has been explored
60 in a variety of settings. Earlier work [e.g., Klare et al., 2012, O’Toole et al., 2012] focuses on
61 single-demographic effects (specifically, race and gender) in pre-deep-learning face detection and
62 recognition. Buolamwini and Gebre [2018] uncovers unequal performance at the phenotypic sub-
63 group level in, specifically, a gender classification task powered by commercial systems. That work,
64 typically referred to as “Gender Shades,” has been and continues to be hugely impactful both within
65 academia and at the industry level. Indeed, Raji and Buolamwini [2019] provide a follow-on analysis,
66 exploring the impact of the Buolamwini and Gebre [2018] paper publicly disclosing performance
67 results, for specific systems, with respect to demographic effects; they find that their named companies
68 (IBM, Microsoft, and Megvii) updated their APIs within a year to address some concerns that were
69 surfaced. Subsequently, the late 2019 update to the NIST FRVT provides evidence that commercial
70 platforms are continuing to focus on performance at the group and subgroup level [Grother et al.,
71 2019]. [Further recent work explores these demographic questions with a focus on Indian election
72 candidates \[Jain and Parsheera, 2021\].](#) We see our benchmark as adding to this literature by, for the
73 first time, addressing both noise and demographic effects on commercial platforms’ face detection
74 offerings.

75 In this work, we focus on *measuring* the impact of noise on a classification task, [like that of Wilber
76 et al. \[2016\]](#); indeed, a core focus of our benchmark is to *quantify* relative drops in performance
77 conditioned on an input datapoint’s membership in a particular group. We view our work as a
78 *benchmark*, that is, it focuses on quantifying and measuring, decidedly not providing a new method
79 to “fix” or otherwise mitigate issues of demographic inequity in a system. Toward that latter point,
80 existing work on “fixing” unfair systems can be split into three (or, arguably, four [Savani et al.,
81 2020]) focus areas: pre-, in-, and post-processing. Pre-processing work largely focuses on dataset
82 curation and preprocessing [e.g., Feldman et al., 2015, Ryu et al., 2018, Quadrianto et al., 2019, Wang
83 and Deng, 2020]. In-processing often constrains the ML training method or optimization algorithm
84 itself [e.g., Zafar et al., 2017b,a, 2019, Donini et al., 2018, Goel et al., 2018, Padala and Gujar, 2020,
85 Agarwal et al., 2018, Wang and Deng, 2020, Martinez et al., 2020, Diana et al., 2020, Lahoti et al.,
86 2020], or focuses explicitly on so-called fair representation learning [e.g., Adeli et al., 2021, Dwork
87 et al., 2012, Zemel et al., 2013, Edwards and Storkey, 2016, Madras et al., 2018, Beutel et al., 2017,
88 Wang et al., 2019]. Post-processing techniques adjust decisioning at inference time to align with
89 quantitative fairness definitions [e.g., Hardt et al., 2016, Wang et al., 2020].

90 **Robustness to noise.** Quantifying, and improving, the robustness to noise of face detection and
91 recognition systems is a decades-old research challenge. Indeed, mature challenges like NIST’s
92 Facial Recognition Vendor Test (FRVT) have tested for robustness since the early 2000s [Phillips
93 et al., 2007]. We direct the reader to a comprehensive introduction to an earlier robustness challenge
94 due to NIST [Phillips et al., 2011]; that work describes many of the specific challenges faced by



Figure 1: Our benchmark consists of 5,066,312 images of the 15 types of algorithmically generated corruptions produced by ImageNet-C. We use data from four datasets (Adience, CCD, MIAP, and UTKFace) and present examples of corruptions from each dataset here.

95 face detection and recognition systems, often grouped into Pose, Illumination, and Expression
 96 (PIE). It is known that commercial systems still suffer from degradation due to noise [e.g., Hosseini
 97 et al., 2017]; none of this work also addresses the intersection of noise with fairness, as we do.
 98 Recently, *adversarial* attacks have been proposed that successfully break commercial face recognition
 99 systems [Shan et al., 2020, Cherepanova et al., 2021]; we note that our focus is on *natural* noise,
 100 as motivated by Hendrycks and Dietterich [2019] by their ImageNet-C benchmark. Literature at
 101 the intersection of adversarial robustness and fairness is nascent and does not address commercial
 102 platforms [e.g., Singh et al., 2020, Nanda et al., 2021]. To our knowledge, our work is the first
 103 systematic benchmark for commercial face detection systems that addresses, comprehensively, noise
 104 and its differential impact on (sub)groups of the population.

105 3 Experimental Description

106 **Datasets and Protocol.** This benchmark uses four datasets to evaluate the robustness of Amazon
 107 AWS and Microsoft Azure’s face detection systems. They are described below.

108 The Open Images Dataset V6 – Extended; More Inclusive Annotations for People (**MIAP**) dataset
 109 [Schumann et al., 2021] was released by Google in May 2021 as an extension of the popular, permissive-
 110 licensed Open Images Dataset specifically designed to improve annotations of humans. For each
 111 image, every human is exhaustively annotated with bounding boxes for the entirety of their person
 112 visible in the image. Each annotation also has perceived gender (Feminine/Masculine/Unknown)
 113 presentation and perceived age (Young, Middle, Old, Unknown) presentation.

114 The Casual Conversations Dataset (**CCD**) [Hazirbas et al., 2021] was released by Facebook in April
 115 2021 under limited license and includes videos of actors. Each actor consented to participate in an
 116 ML dataset and provided their self-identification of age and gender (coded as Female, Male, and
 117 Other), each actor’s skin type was rated on the Fitzpatrick scale [Fitzpatrick, 1988], and each video
 118 was rated for its ambient light quality. For our benchmark, we extracted one frame from each video.

119 The **Adience** dataset [Eidinger et al., 2014] under a CC license, includes cropped images of faces
 120 from images “in the wild”. Each cropped image contains only one primary, centered face, and each
 121 face is annotated by an external evaluator for age and gender (Female/Male). The ages are reported
 122 as member of 8 age range buckets: 0-2; 3-7; 8-14; 15-24; 25-35; 36-45; 46-59; 60+.

123 Finally, the **UTKFace** dataset [Zhang et al., 2017] under a non-commercial license, contains images
 124 with one primary subject and were annotated for age (continuous), gender (Female/Male), and
 125 ethnicity (White/Black/Asian/Indian/Others) by an algorithm, then checked by human annotators.

126 For each of the datasets, we randomly selected a subset of images for our evaluation in order to cap
 127 the number of images from each intersectional identity at 1,500 as an attempt to reduce the effect of
 128 highly imbalanced datasets. We include a total of 66,662 images with 14,919 images from Adience;
 129 21,444 images from CCD; 8,194 images from MIAP; and 22,105 images from UTKFace. The full
 130 breakdown of totals of images from each group can be found in Section A.1.

131 Each image was corrupted a total of 75 times, per the ImageNet-C protocol with the main 15
 132 corruptions each with 5 severity levels. Examples of these corruptions can be seen in Figure 1. This
 133 resulted in a total of 5,066,312 images (including the original clean ones) which were each passed
 134 through the AWS and Azure face analysis systems. A detailed description of which API settings were
 135 selected can be found in Appendix C. The API calls were conducted between 19 May and 29 May
 136 2021. Images were processed and stored within AWS’s cloud using S3 and EC2. The total cost of the
 137 experiments was \$9,887.17 and a breakdown of costs can be found in Appendix D.

138 **Evaluation Metrics.** Given that we aim is to study how corruptions to an image alter the commer-
 139 cial interpretation of that image, we evaluate the error of the face systems. Additionally, none of the
 140 chosen datasets have ground truth face bounding boxes. Therefore, we can use the response from the
 141 clean image as a ground truth of sorts. Specifically, we take as ground truth the number of faces in an
 142 clean image and compare that to the number of faces detected in a corrupted image.

143 Our main metric is the relative error in the number of faces a system detects after corruption; [this](#)
 144 [metric has been used in other facial processing benchmarks \[Jain and Parsheera, 2021\]. Measuring](#)
 145 [error in this way is in some sense incongruous with the object detection nature of the APIs. However,](#)
 146 [none of the data in our datasets have bounding boxes for each face. This means that we cannot](#)
 147 [calculate precision metrics as one would usually do with other detection tasks. To overcome this,](#)
 148 [we hand annotated bounding boxes for each face in 772 random images from the dataset. We then](#)
 149 [calculated per-image precision scores \(with an IoU of 0.5\) and per-image relative error in face counts](#)
 150 [and we find a Pearson’s correlation of 0.91 \(with \$p < 0.001\$ \). This high correlation indicates that the](#)
 151 [proxy is sufficient to be used in this benchmark in the absence of fully annotated bounding boxes.](#)

152 This error is calculated for each image. The way in which this works is that we first pass every clean,
 153 uncorrupted image through the commercial system’s API. Then, we measure the number of detected
 154 faces, i.e., length of the system’s response, and treat this number as the ground truth. Subsequently,
 155 we compare the [number of detected faces](#) for a corrupted version of that image. If the two [face counts](#)
 156 are not the same, then we call that an error. We refer to this as the *relative corruption error*. For each
 157 clean image, i , from dataset d , and each corruption c which produces a corrupted image $\hat{i}_{c,s}$ with
 158 severity s , we compute the relative corruption error for system r as

$$rCE_{c,s}^{d,r}(\hat{i}_{c,s}) := \begin{cases} 1, & \text{if } l_r(i) \neq l_r(\hat{i}_{c,s}) \\ 0, & \text{if } l_r(i) = l_r(\hat{i}_{c,s}) \end{cases}$$

159 where l_r computes the number of detected faces, i.e., length of the response, from face detection
 160 system r when given an image. Often the super- and subscripts are omitted when they are obvious
 161 from context.

162 Our main metric, relative error, aligns with that of the ImageNet-C benchmark. We report mean
 163 relative corruption error ($mrCE$) defined as taking the average of rCE across some relative set
 164 of categories. In our experiments, depending on the context, we might have any of the following
 165 categories: face systems, datasets, corruptions, severities, age presentation, gender presentation,
 166 Fitzpatrick rating, and ambient lighting. For example, we might report the relative mean corruption
 167 error when averaging across demographic groups; the mean corruption error for Azure on the UTK
 168 dataset for each age group a is $mrCE_a = \frac{1}{15} \frac{1}{5} \sum_{c,s} rCE_{c,s,a}^{UTK,Azure}$. The subscripts on $mrCE$
 169 will be omitted when it is obvious what their value is in whatever context they are presented.

170 Finally, we will also investigate the significance of whether the $mrCE$ for two groups are equal. For
 171 example, our first question is whether the two commercial systems (AWS and Azure) have comparable
 172 $mrCE$ overall. To do this, we will report the raw $mrCE$; these frequency or empiric probability
 173 statistics offer much insight into the likelihood of error. But we also indicate the statistical significance
 174 at $\alpha = 0.05$ determined by logistic regressions for the appropriate variables and interactions. For
 175 each claim of significance, regression tables can be found in the appendix. Accordingly, we discuss
 176 the odds or odds ratio of relevant features. See Appendix B for a detailed example. **Finally, each**

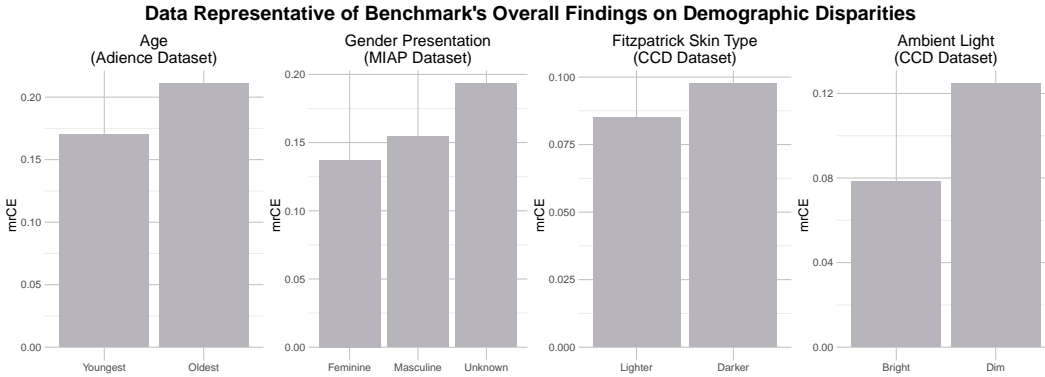


Figure 2: *There are disparities in all of the demographics included in this study; we show representative evidence for each demographic on different datasets. On the left, we see (using Adience as an exemplar) that the oldest two age groups are roughly 25% more error prone than the youngest two groups. Using MIAP as an exemplar, masculine presenting subjects are 20% more error prone than feminine. On the CCD dataset, we find that individuals with Fitzpatrick scales IV-VI have a roughly 25% higher chance of error than lighter skinned individuals. Finally, dimly lit individuals are 60% more likely to have errors.*

177 claim we make for an individual dataset or service is backed up with statistical rigor through the
 178 logistic regressions. Each claim we make across datasets is done by looking at the trends in each
 179 dataset and are inherently qualitative.

180 **What is not included in this study.** There are three main things that this benchmark does not
 181 address. First, we do not examine cause and effect. We report inferential statistics without discussion
 182 of what generates them. Second, we only examine the types of algorithmically generated natural
 183 noise present in the 15 corruptions. We speak narrowly about robustness to these corruptions or
 184 perturbations. We explicitly do not study or measure robustness to other types of changes to images,
 185 for instance adversarial noise, camera dimensions, etc. Finally, we do not investigate algorithmic
 186 training. We do not assume any knowledge of how the commercial system was developed or what
 187 training procedure or data were used.

188 **Social Context.** The central analysis of this benchmark relies on socially constructed concepts
 189 of gender presentation and the related concepts of race and age. While this benchmark analyzes
 190 phenotypal versions of these from metadata on ML datasets, it would be wrong to interpret our
 191 findings absent a social lens of what these demographic groups mean inside a society. We guide the
 192 reader to Benthall and Haynes [2019] and Hanna et al. [2020] for a look at these concepts for race in
 193 machine learning, and Hamidi et al. [2018] and Keyes [2018] for similar looks at gender.

194 4 Benchmark Results

195 We now report the main results of our benchmark, a synopsis of which is in Figure 2. Overall, we
 196 find that photos of individuals who are *older*, *masculine presenting*, *darker skinned*, or are *dimly lit*
 197 are more susceptible to errors than their counterparts. We come to these qualitative conclusions by
 198 quantitatively examining the trends of each dataset for each demographic. All four datasets have
 199 age and gender labels. We see the bias against older individuals across all datasets. The bias against
 200 masculine presenting individuals is present in all datasets except UTKFace (which shows no bias).
 201 Skin type and lighting labels are only present in one dataset, CCD.

202 Below is a more detailed analysis with additional supporting tables and figures in the Appendix.

203 4.1 System Performance

204 Overall, AWS has fewer errors than Azure on corrupted data though the magnitude of the difference
 205 is small. The *mrCE* for AWS is 12.298% whereas Azure is 12.338%, or 3% higher, but this is a
 206 Simpson’s Paradox because when we look at each dataset, we see further nuance.

207 We plot the CTRs for each dataset and service in Figure 3; the difference between services is
 208 statistically significant for each dataset. For the Adience and MIAP datasets, Azure performs better

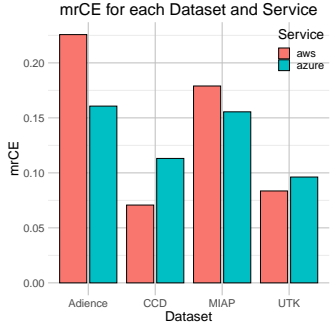


Figure 3: Observe that AWS is more robust on CCD and UTK and Azure is more robust on Adience and MIAP.

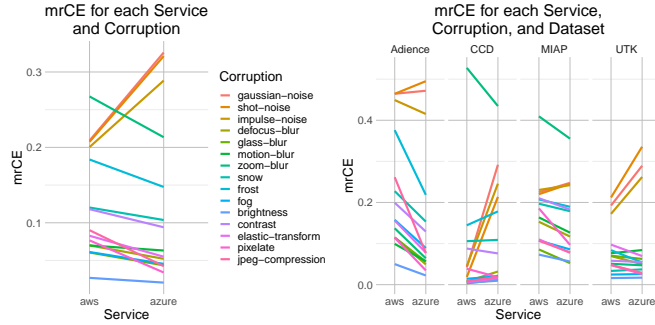


Figure 4: A comparison of $mrCE$ for each commercial system and dataset where each line represents one of the 15 types of corruptions. (Left) depicts the robustness across all datasets whereas (right) depicts this for each dataset separately.

209 than AWS. On Adience, Azure’s $mrCE$ is 16.1% whereas AWS has $mrCE$ of 22.6%. The magnitude
 210 is less on MIAP; Azure has 15.6% and AWS has 17.9%.

211 Conversely, on the CCD and UTK dataset, Azure outperforms AWS. For the CCD dataset, Azure
 212 performs 60% worse than AWS (AWS $mrCE$ of 7.1% compared to Azure’s 11.3%). The magnitude
 213 is less on UTKFace; AWS has 8.4% whereas Azure has 9.6%.

214 4.2 Noise corruptions are the most difficult

215 Recall that the ImageNet-C corruptions are broken into four different types: noise, blur, weather,
 216 and digital corruptions. We observe that the noise corruptions prove to be some of the most difficult
 217 corruptions for the commercial systems to handle. From Figure 4, we observe that in the AWS system,
 218 the three noise corruptions have the the second, third, and fourth most difficult corruptions (behind
 219 zoom blur). However, they are markedly the most difficult corruptions for Azure to handle. On the
 220 otherhand, Azure outperforms AWS on every other corruption. The difficulty of the noise corruptions
 221 echos that documented in the ImageNet-C experiments, though the comparative magnitude of the
 222 difficulty for these systems is significantly higher than what is previously documented.

223 When we examine the differences in the performance for each corruption across the different datasets,
 224 we see a continuation of the theme that the noise corruptions have relatively high $mrCE$. In every
 225 instance except one, Azure performs worse on the noise corruptions than AWS. For both commercial
 226 systems on Adience, the $mrCE$ values for the noise corruptions are above 40%. However, Azure
 227 preforms better than AWS on all other corruptions on the Adience Dataset.

228 The zoom blur corruption proves particularly difficult on the CCD and MIAP datasets, though Azure
 229 is significantly better than AWS (CCD: 52.7% for AWS and 43.5% for Azure; MIAP: 41.0% for
 230 AWS and 35.5% for Azure). We also note that all corruptions for all datasets and commercial systems
 231 are significantly differently from zero.

232 4.2.1 Comparison to ImageNet-C results

233 Even though Hendrycks and Dietterich [2019] worked with the ImageNet dataset, we compare the
 234 findings from their paper to our experiments. We recreate Figure 3 from their paper with more current
 235 results for recent models since their paper was published, as well as the addition of our findings for
 236 AWS and Azure’s face detection on our data; see Figure 8. This figure reproduces their metric, mean
 237 corruption error and relative mean corruption error. These differ from our metrics as they are defined
 238 as the raw error for each corruption, but normalized against the performance of AlexNet from the
 239 original paper. This is done so as to compare different models more fairly. The figure also shows the
 240 relative mean corruption error which is the difference between the raw error for each corruption and
 241 the raw error for the clean data. From this figure, we can conclude that our results are very highly
 242 in-line with the predictions from the previous data. This indicates that, even with highly accurate
 243 models, accuracy is a strong predictor of robustness.

244 We also examined the corruption-specific differences between our findings (with face data) and that
 245 of the original paper (with ImageNet data). We find that while facial datasets are most susceptible to

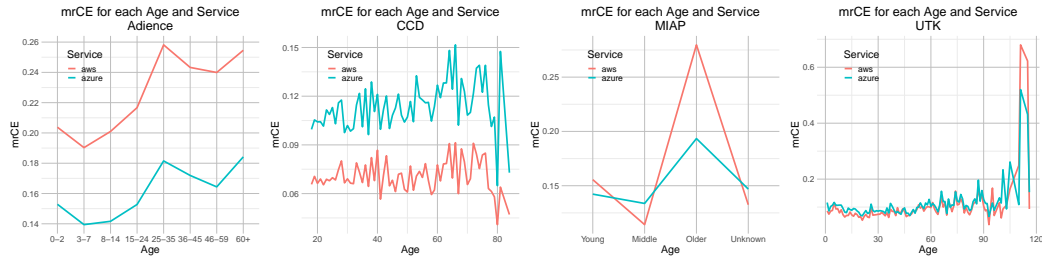


Figure 5: Each figure depicts the $mrCE$ across ages. Each line depicts a commercial system (AWS is above Azure for Adience and MIAP). Age is a categorical variable for Adience and MIAP but a numeric for CCD and UTKFace. Observe the general trend of increased errors for older individuals.

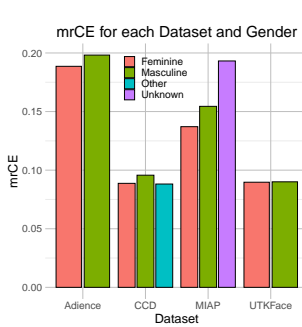


Figure 6: Observe that on all datasets, except for UTKFace, feminine presenting individuals are more robust than masculine.

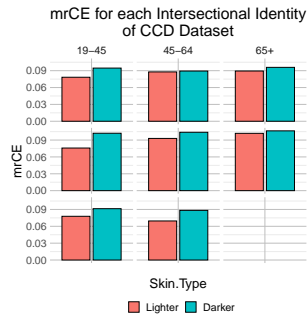


Figure 7: In all intersectional identities, except for 45-64 females, darker skinned individuals are less robust than those who are lighter skinned.

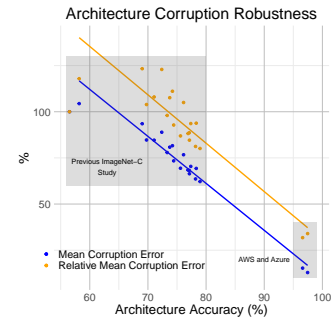


Figure 8: Recreation of Figure 3 from Hendrycks and Dietterich [2019] with new results since their paper and the addition of our findings.

246 noise corruptions, zoom blur, weather, etc, the ImageNet datasets are generally uniformly susceptible
 247 to corruptions with blurs and digital corruptions being the most difficult for them. This indicates
 248 that the face data have qualitative differences in their robustness susceptibility, indicating a need for
 249 further study.

250 4.3 Errors increase on older subjects

251 We observe a significant impact of age on $mrCE$. See Figure 5. In every dataset and every
 252 commercial system, we see that older subjects have significantly higher error rates. Recall that all
 253 four datasets have age metadata. Adience and MIAP have such data in groups. CCD and UTKFace
 254 have age data as a continuous variable.

255 On the Adience dataset, there is an interesting behavior where the second and third youngest age
 256 groups have the best performance with increases for younger and older age groups. There is then
 257 a spike in errors in the 25-35 age group which falls off slightly for the 36-59 groups and finally
 258 increases again for the oldest 60+ group. These two maximal groups have nearly 1:4 odds of error.
 259 This is compared to the youngest group which has 30% better odds (3:15).

260 For the MIAP dataset, the age disparity is very pronounced. Like the Adience dataset, we see a
 261 decrease in the likelihood of error moving from the youngest to the middle ages. However, we see a
 262 very large increase for the Oldest individuals. In AWS for instance, we see a 145% increase in error.

263 The CCD and UTKFace datasets have numeric age. Analyzing the regressions indicates that for every
 264 increase of 10 years, there is a 2.3% increase in the likelihood of error on the CCD data and 2.7%
 265 increase for UTKFace data. In Appendix E.4, we explore the interaction of Age and the corruptions.

266 4.4 Masculine presenting individuals have more errors than feminine presenting

267 Across all datasets except UTKFace, we find that feminine presenting individuals have lower errors
 268 than masculine presenting individuals. See Figure 6. On Adience, feminine individuals have 18.8%
 269 $mrCE$ whereas masculine have 19.8%. On CCD, the $mrCE$ s are 8.9% and 9.6% respectively. On

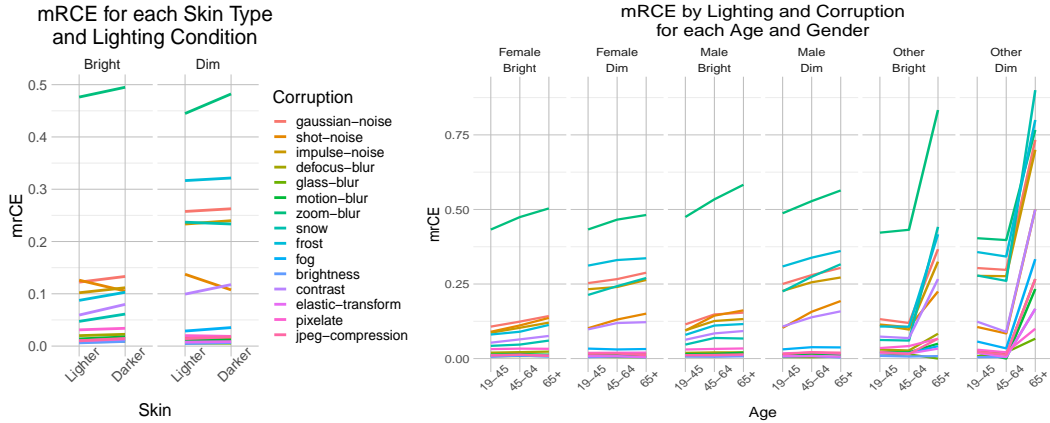


Figure 9: (Left) $mrCE$ is plotted for each corruption by the intersection of lighting condition and skin type. (Right) the same is plotted by the intersection of age, gender, and lighting. Observe that for both skin types, all genders, and all ages, the dimly lit environment increases the error rates. Motion blur is the least robust corruption with frost, the three noises, and snow being the next worst across most intersectional identities.

270 the MIAP dataset, the $mrCE$ values are 13.7% and 15.4% respectively. On the UTKFace, both
 271 gender presentations have around 9.0% $mrCE$ (non statistically significant difference).

272 Stepping outside the gender binary, we have two insights into this from these data. In the CCD
 273 dataset, the subjects were asked to self-identify their gender. Two individuals selected Other and 62
 274 others did not provide a response. Those two who chose outside the gender binary have a $mrCE$
 275 of 4.9%. When we include those individuals without gender labels, their $mrCE$ is 8.8% and not
 276 significantly different from the feminine presenting individuals.

277 The other insight comes from the MIAP dataset where subjects were rated on their perceived
 278 gender presentation by crowdworkers; options were “Predominantly Feminine”, “Predominantly
 279 Masculine”, and “Unknown”. For those “Unknown”, the overall $mrCE$ is 19.3%. The creators of the
 280 dataset automatically set the gender presentation of those with an age presentation of “Young” to be
 281 “Unknown”. The $mrCE$ of those annotations which aren’t “Young” and have an “Unknown” gender
 282 presentation raises to 19.9%. One factor that might contribute to this phenomenon is that individuals
 283 with an “Unknown” gender presentation might have faces that are occluded or are small in the image.
 284 Further work should be done to explore the causes of his discrepancy. In Appendix E.3, we explore
 285 the interaction of Gender and the corruptions.

286 4.5 Dark skinned subjects have more errors across age and gender identities

287 We analyze data from the CCD dataset which has ratings for each subject on the Fitzpatrick scale.
 288 As is customary in analyzing these ratings, we split the six Fitzpatrick values into two: Lighter (for
 289 ratings I-III) and Darker for ratings (IV-VI). The main intersectional results are reported in Figure 7.

290 The overall $mrCE$ for lighter and darker skin types are 8.5% and 9.7% respectively, a 15% increase
 291 for the darker skin type. We also see a similar trend in the intersectional identities available in the
 292 CCD metadata (age, gender, and skin type). We see that in every identity (except for 45-64 year old
 293 and Feminine) the darker skin type has statistically significant higher error rates. This difference
 294 is particularly stark in 19-45 year old, masculine subjects. We see a 35% increase in errors for the
 295 darker skin type subjects in this identity compared to those with lighter skin types. For every 20
 296 errors on a light skinned, masculine presenting individual between 18 and 45, there are 27 errors for
 297 dark skinned individuals of the same category.

298 4.6 Dim lighting conditions has the most severe impact on errors

299 Using lighting condition information from the CCD dataset, we observe the $mrCE$ is substantially
 300 higher in dimly lit environments: 12.5% compared to 7.8% in bright environments. See Figure 9.

301 Across the board, we generally see that the disparity in demographic groups decreases between bright
 302 and dimly lit environments. For example, the odds ratio between dark and light skinned subjects is
 303 1.09 for bright environments, but decreases to 1.03 for dim environments. This is true for age groups

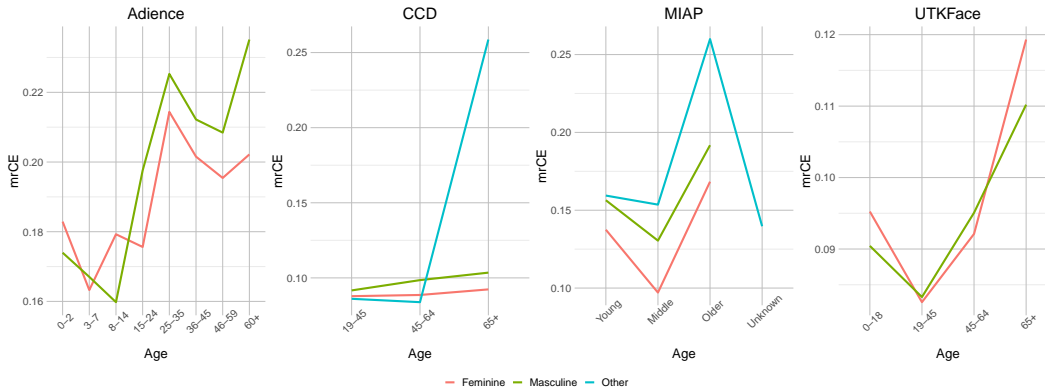


Figure 10: For each dataset, the $mrCE$ is plotted across age groups. Each gender is represented and indicates how gender disparities change across the age groups.

304 (e.g., odds ratios 1.183 (bright) vs 1.127 (dim) for 45-64 compared to 19-45; 1.138 (bright) vs 1.060
 305 (dim) for Males compared to Females). This is not true for individuals with gender identities as Other
 306 or omitted – the disparity increases (1.104 (bright) vs 1.173 (dim) with Females as the reference).

307 In Figure 9 we observe the lighting differences for different intersectional identities across corruptions.
 308 We continue to see zoom blur as the most challenging corruption. Interestingly, the noise and some
 309 weather corruptions have a large increase in their errors in dimly lit environments across intersectional
 310 identities whereas many of the other corruptions do not.

311 4.7 Older subjects have higher gender error disparities

312 We plot in Figure 10 the $mrCE$ for each dataset across age with each gender group plotted separately.
 313 From this, we can note that on the CCD and MIAP dataset, the masculine presenting group is always
 314 less robust than the feminine. On the CCD dataset, the disparity between the two groups increases
 315 as the age increases (odds ratio of 1.048 for 19-45 raises to 1.135 for 65+). On the MIAP dataset,
 316 the odds ratio is greatest between masculine and feminine for the middle age group (1.395). The
 317 disparities between the ages also increases from feminine to masculine to unknown gender identities.

318 On the Adience and UTKFace datasets, we see that the feminine presenting individuals sometime
 319 have higher error rates than masculine presenting subjects. Notably, the most disparate errors in
 320 genders on these datasets occurs at the oldest categories, following the trend from the other datasets.

321 5 Gender and Age Estimation Analysis

322 We briefly overview results from evaluating AWS’s age and gender estimation commercial systems.
 323 The detection model we evaluated for Azure does not provide age and gender estimates. Further
 324 analysis can be found in Appendices F and G.

325 5.1 Gender estimation is at least twice as susceptible to corruptions as face detection

326 The use of automated gender estimates in ML is a controversial topic. Trans and gender queer
 327 individuals are often ignored in ML research, though there is a growing body of research that aims
 328 to use these technologies in an assistive way as well [e.g., Ahmed, 2019, Chong et al., 2021]. To
 329 evaluate gender estimation, we only use CCD as the subjects of these photos voluntarily identified
 330 their gender. We omit from the analysis any individual who either did not choose to give their gender
 331 or fall outside the gender binary because AWS only estimates Male and Female.

332 AWS misgenders 9.1% of the clean images but 21.6% of the corrupted images. Every corruption
 333 performs worse on gender estimation than $mrCE$. Two corruptions (elastic transform and glass
 334 blur) do not have statistically different errors from the clean images. All the others do, with the most
 335 significant being zoom blur, Gaussian noise, impulse noise, snow, frost, shot noise, and contrast.
 336 Zoom blur’s probability of error is 61% and Gaussian noise is 32%. This compares to $mrCE$ values
 337 of 43% and 29% respectively. See Appendix F for further analysis.

338 **5.2 Corrupted images error in their age predictions by 40% more than clean images**

339 To estimate Age, AWS returns an upper and lower age estimation. Following their own guidelines on
340 face detection,¹ we use the mid-point of these numbers as a approximate estimate. On average, the
341 estimation is 8.3 years away from the actual age of the subject for corrupted data, this compares to
342 5.9 years away for clean data. See Appendix G for further analysis.

343 **6 Conclusion**

344 This benchmark has evaluated two leading commercial facial detection and analysis systems for their
345 robustness against common natural noise corruptions. Using the 15 ImageNet-C corruptions, we
346 measured the relative mean corruption error as measured by comparing the number of faces detected
347 in a clean and corrupted image. We used four academic datasets which included demographic detail.
348 Adience, MIAP, and UTKFace have perceived age and gender metadata. CCD has subject provided
349 age and gender responses as well as external ratings of skin type and ambient lighting conditions.

350 We observed through our analysis that there are significant demographic disparities in the likelihood
351 of error on corrupted data. We found that older individuals, masculine presenting individuals, those
352 with darker skin types, or in photos with dim ambient light all have higher errors ranging from
353 20-60%. We also investigated questions of intersectional identities finding that darker males have
354 the highest corruption errors. As for age and gender estimation, corruptions have a significant
355 and sizeable impact on the system's performance; gender estimation is more than twice as bad on
356 corrupted images as it is on clean images; age estimation is 40% worse on corrupted images.

357 Future work could explore other metrics for evaluating face detection systems when ground truth
358 bounding boxes are not present. While we considered the length of response on clean images to be
359 ground truth, it could be viable to treat the clean image's bounding boxes as ground truth and measure
360 deviations therefrom when considering questions of robustness. Of course, this would require a
361 transition to detection-based metrics like precision, recall, and F -measure.

362 We do not explore questions of causation in this benchmark. We do not have enough different
363 datasets or commercial systems to probe this question through regressions or mixed effects modeling.
364 We do note that there is work that examines causation questions with such methods like that of
365 [Best-Rowden and Jain, 2017] and [Cook et al., 2019]. With additional data and under similar
366 benchmarking protocols, one could start to examine this question. However, the black-box nature of
367 commercial systems presents unique challenges to this endeavor.

¹<https://docs.aws.amazon.com/rekognition/latest/dg/guidance-face-attributes.html>

368 **References**

- 369 E. Adeli, Q. Zhao, A. Pfefferbaum, E. V. Sullivan, L. Fei-Fei, J. C. Niebles, and K. M. Pohl.
370 Representation learning with statistical independence to mitigate bias. In *Proceedings of the*
371 *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2513–2523, 2021.
- 372 A. Agarwal, A. Beygelzimer, M. Dudik, J. Langford, and H. Wallach. A reductions approach to
373 fair classification. In *Proceedings of the 35th International Conference on Machine Learning*,
374 volume 80, pages 60–69, 2018. URL <http://proceedings.mlr.press/v80/agarwal18a.html>.
375
- 376 A. A. Ahmed. Bridging social critique and design: Building a health informatics tool for transgender
377 voice. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*,
378 pages 1–4, 2019.
- 379 S. Barocas, M. Hardt, and A. Narayanan. *Fairness and Machine Learning*. fairmlbook.org, 2019.
380 <http://www.fairmlbook.org>.
- 381 S. Benthall and B. D. Haynes. Racial categories in machine learning. In *Proceedings of the conference*
382 *on fairness, accountability, and transparency*, pages 289–298, 2019.
- 383 L. Best-Rowden and A. K. Jain. Longitudinal study of automatic face recognition. *IEEE transactions*
384 *on pattern analysis and machine intelligence*, 40(1):148–162, 2017.
- 385 A. Beutel, J. Chen, Z. Zhao, and E. H. Chi. Data decisions and theoretical implications when
386 adversarially learning fair representations. *arXiv preprint arXiv:1707.00075*, 2017.
- 387 J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial
388 gender classification. In *Proceedings of the 1st Conference on Fairness, Accountability and*
389 *Transparency*, volume 81, pages 77–91, 2018. URL <http://proceedings.mlr.press/v81/buolamwini18a.html>.
390
- 391 V. Cherepanova, M. Goldblum, H. Foley, S. Duan, J. P. Dickerson, G. Taylor, and T. Goldstein.
392 Lowkey: leveraging adversarial attacks to protect social media users from facial recognition. In
393 *International Conference on Learning Representations (ICLR)*, 2021.
- 394 T. Chong, N. Maudet, K. Harima, and T. Igarashi. Exploring a makeup support system for transgender
395 passing based on automatic gender recognition. In *Proceedings of the 2021 CHI Conference on*
396 *Human Factors in Computing Systems*, pages 1–13, 2021.
- 397 A. Chouldechova and A. Roth. The frontiers of fairness in machine learning. *arXiv preprint*
398 *arXiv:1810.08810*, 2018.
- 399 C. M. Cook, J. J. Howard, Y. B. Sirotin, J. L. Tipton, and A. R. Vemury. Demographic effects in
400 facial recognition and their dependence on image acquisition: An evaluation of eleven commercial
401 systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 1(1):32–41, 2019.
- 402 W. Derringer. A surveillance net blankets china’s cities, giving police vast powers. *The New York*
403 *Times*, Dec. 17 2019. URL <https://www.nytimes.com/2019/12/17/technology/china-surveillance.html>.
404
- 405 E. Diana, W. Gill, M. Kearns, K. Kenthapadi, and A. Roth. Convergent algorithms for (relaxed)
406 minimax fairness. *arXiv preprint arXiv:2011.03108*, 2020.
- 407 M. Donini, L. Oneto, S. Ben-David, J. Shawe-Taylor, and M. Pontil. Empirical risk minimization
408 under fairness constraints. In *Proceedings of the 32nd International Conference on Neural*
409 *Information Processing Systems, NIPS’18*, page 2796–2806, 2018.
- 410 C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel. Fairness through awareness. In *Proceed-*
411 *ings of the 3rd Innovations in Theoretical Computer Science Conference, ITCS ’12*, page 214–226,
412 New York, NY, USA, 2012. Association for Computing Machinery. ISBN 9781450311151. doi:
413 10.1145/2090236.2090255. URL <https://doi.org/10.1145/2090236.2090255>.

- 414 H. Edwards and A. J. Storkey. Censoring representations with an adversary. In *4th International*
415 *Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016,*
416 *Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1511.05897>.
- 417 E. Eidinge, R. Enbar, and T. Hassner. Age and gender estimation of unfiltered faces. *IEEE*
418 *Transactions on Information Forensics and Security*, 9(12):2170–2179, 2014.
- 419 H. El Khiyari and H. Wechsler. Face verification subject to varying (age, ethnicity, and gender)
420 demographics using deep learning. *Journal of Biometrics and Biostatistics*, 7(323):11, 2016.
- 421 M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian. Certifying and
422 removing disparate impact. In *Knowledge Discovery and Data Mining*, pages 259–268, 2015.
- 423 T. B. Fitzpatrick. The validity and practicality of sun-reactive skin types i through vi. *Archives of*
424 *dermatology*, 124(6):869–871, 1988.
- 425 C. Garvie. *The perpetual line-up: Unregulated police face recognition in America*. Georgetown Law,
426 Center on Privacy & Technology, 2016.
- 427 N. Goel, M. Yaghini, and B. Faltings. Non-discriminatory machine learning through convex fairness
428 criteria. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 2018. URL
429 <https://ojs.aaai.org/index.php/AAAI/article/view/11662>.
- 430 Google. How google uses pattern recognition to make sense of images. [https://policies.](https://policies.google.com/technologies/pattern-recognition?hl=en-US)
431 [google.com/technologies/pattern-recognition?hl=en-US](https://policies.google.com/technologies/pattern-recognition?hl=en-US), 2021. Accessed: 2021-06-
432 07.
- 433 P. Grother, M. Ngan, and K. Hanaoka. *Face Recognition Vendor Test (FVRT): Part 3, Demographic*
434 *Effects*. National Institute of Standards and Technology, 2019.
- 435 D. Gutman. King County Council bans use of facial recognition technology by Sheriff’s Office, other
436 agencies. *The Seattle Times*, June 2021. URL [https://www.seattletimes.com/seattle-](https://www.seattletimes.com/seattle-news/politics/king-county-council-bans-use-of-facial-recognition-technology-by-sheriffs-office-other-agencies/)
437 [news/politics/king-county-council-bans-use-of-facial-recognition-](https://www.seattletimes.com/seattle-news/politics/king-county-council-bans-use-of-facial-recognition-technology-by-sheriffs-office-other-agencies/)
438 [technology-by-sheriffs-office-other-agencies/](https://www.seattletimes.com/seattle-news/politics/king-county-council-bans-use-of-facial-recognition-technology-by-sheriffs-office-other-agencies/).
- 439 F. Hamidi, M. K. Scheuerman, and S. M. Branham. Gender recognition or gender reductionism?
440 the social implications of embedded gender recognition systems. In *Proceedings of the 2018 chi*
441 *conference on human factors in computing systems*, pages 1–13, 2018.
- 442 A. Hanna, E. Denton, A. Smart, and J. Smith-Loud. Towards a critical race methodology in
443 algorithmic fairness. In *Proceedings of the 2020 conference on fairness, accountability, and*
444 *transparency*, pages 501–512, 2020.
- 445 M. Hardt, E. Price, E. Price, and N. Srebro. Equality of opportunity in super-
446 vised learning. In *Advances in Neural Information Processing Systems*, volume 29,
447 pages 3315–3323, 2016. URL [https://proceedings.neurips.cc/paper/2016/file/](https://proceedings.neurips.cc/paper/2016/file/9d2682367c3935defcb1f9e247a97c0d-Paper.pdf)
448 [9d2682367c3935defcb1f9e247a97c0d-Paper.pdf](https://proceedings.neurips.cc/paper/2016/file/9d2682367c3935defcb1f9e247a97c0d-Paper.pdf).
- 449 W. Hartzog. The secretive company that might end privacy as we know it. *The New York Times*, Jan. 18
450 2020. URL [https://www.nytimes.com/2020/01/18/technology/clearview-privacy-](https://www.nytimes.com/2020/01/18/technology/clearview-privacy-facial-recognition.html)
451 [facial-recognition.html](https://www.nytimes.com/2020/01/18/technology/clearview-privacy-facial-recognition.html).
- 452 C. Hazirbas, J. Bitton, B. Dolhansky, J. Pan, A. Gordo, and C. C. Ferrer. Towards measuring fairness
453 in ai: the casual conversations dataset. *arXiv preprint arXiv:2104.02821*, 2021.
- 454 D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions
455 and perturbations. 2019.
- 456 H. Hosseini, B. Xiao, and R. Poovendran. Google’s cloud vision API is not robust to noise. In
457 *2017 16th IEEE international conference on machine learning and applications (ICMLA)*, pages
458 101–105. IEEE, 2017.
- 459 G. Jain and S. Parsheera. 1.4 billion missing pieces? auditing the accuracy of facial processing tools
460 on indian faces. *First Workshop on Ethical Considerations in Creative applications of Computer*
461 *Vision*, 2021.

- 462 S. Kantayya. Coded bias, 2020. Feature-length documentary.
- 463 O. Keyes. The misgendering machines: Trans/hci implications of automatic gender recognition.
464 *Proceedings of the ACM on human-computer interaction*, 2(CSCW):1–22, 2018.
- 465 B. F. Klare, M. J. Burge, J. C. Klontz, R. W. V. Bruegge, and A. K. Jain. Face recognition performance:
466 Role of demographic information. *IEEE Transactions on Information Forensics and Security*, 7(6):
467 1789–1801, 2012.
- 468 P. Lahoti, A. Beutel, J. Chen, K. Lee, F. Prost, N. Thain, X. Wang, and E. H. Chi. Fairness without
469 demographics through adversarially reweighted learning. *arXiv preprint arXiv:2006.13114*, 2020.
- 470 S. Lohr. Facial recognition is accurate, if you’re a white guy. *New York Times*, 9, 2018.
- 471 D. Madras, E. Creager, T. Pitassi, and R. S. Zemel. Learning adversarially fair and transferable
472 representations. In *Proceedings of the 35th International Conference on Machine Learning, ICML*
473 *2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of*
474 *Machine Learning Research*, pages 3381–3390. PMLR, 2018. URL [http://proceedings.mlr.
475 press/v80/madras18a.html](http://proceedings.mlr.press/v80/madras18a.html).
- 476 J. Marson and B. Forrest. Armed low-cost drones, made by turkey, reshape battlefields and geopolitics.
477 *The Wall Street Journal*, Jun 2021. URL [https://www.wsj.com/articles/armed-low-cost-
478 drones-made-by-turkey-reshape-battlefields-and-geopolitics-11622727370](https://www.wsj.com/articles/armed-low-cost-drones-made-by-turkey-reshape-battlefields-and-geopolitics-11622727370).
- 479 N. Martinez, M. Bertran, and G. Sapiro. Minimax pareto fairness: A multi objective perspective.
480 In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages
481 6755–6764, 2020. URL <http://proceedings.mlr.press/v119/martinez20a.html>.
- 482 V. Nanda, S. Dooley, S. Singla, S. Feizi, and J. P. Dickerson. Fairness through robustness: Investi-
483 gating robustness disparity in deep learning. In *Proceedings of the 2021 ACM Conference on*
484 *Fairness, Accountability, and Transparency*, pages 466–477, 2021.
- 485 A. J. O’Toole, P. J. Phillips, X. An, and J. Dunlop. Demographic effects on estimates of automatic
486 face recognition performance. *Image and Vision Computing*, 30(3):169–176, 2012.
- 487 M. Padala and S. Gujar. Fnncc: Achieving fairness through neural networks. In *Proceedings*
488 *of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages
489 2277–2283. International Joint Conferences on Artificial Intelligence Organization, 7 2020. doi:
490 10.24963/ijcai.2020/315. URL <https://doi.org/10.24963/ijcai.2020/315>.
- 491 P. J. Phillips, W. T. Scruggs, A. J. O’Toole, P. J. Flynn, K. W. Bowyer, C. L. Schott, and M. Sharpe.
492 Frvt 2006 and ice 2006 large-scale results. *National Institute of Standards and Technology, NISTIR*,
493 7408(1):1, 2007.
- 494 P. J. Phillips, J. R. Beveridge, B. A. Draper, G. Givens, A. J. O’Toole, D. S. Bolme, J. Dunlop,
495 Y. M. Lui, H. Sahibzada, and S. Weimer. An introduction to the good, the bad, & the ugly face
496 recognition challenge problem. In *2011 IEEE International Conference on Automatic Face &*
497 *Gesture Recognition (FG)*, pages 346–353. IEEE, 2011.
- 498 N. Quadrianto, V. Sharmanska, and O. Thomas. Discovering fair representations in
499 the data domain. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 8227–8236. Com-
500 puter Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00842. URL [http:
501 //openaccess.thecvf.com/content_CVPR_2019/html/Quadrianto_Discovering_
502 Fair_Representations_in_the_Data_Domain_CVPR_2019_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Quadrianto_Discovering_Fair_Representations_in_the_Data_Domain_CVPR_2019_paper.html).
- 504 I. D. Raji and J. Buolamwini. Actionable auditing: Investigating the impact of publicly naming
505 biased performance results of commercial ai products. In *Proceedings of the 2019 AAAI/ACM*
506 *Conference on AI, Ethics, and Society*, pages 429–435, 2019.
- 507 H. J. Ryu, H. Adam, and M. Mitchell. Inclusivefacenet: Improving face attribute detection with race
508 and gender diversity. *arXiv preprint arXiv:1712.00193*, 2018.

- 509 Y. Savani, C. White, and N. S. Govindarajulu. Intra-processing methods for debiasing neural networks.
510 In *Proceedings of Advances in Neural Information Processing Systems*, 2020.
- 511 C. Schumann, C. R. Pantofaru, S. Ricco, U. Prabhu, and V. Ferrari. A step toward more inclusive
512 people annotations for fairness. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and*
513 *Society*, 2021.
- 514 S. Shan, E. Wenger, J. Zhang, H. Li, H. Zheng, and B. Y. Zhao. Fawkes: Protecting privacy against
515 unauthorized deep learning models. In *29th {USENIX} Security Symposium ({USENIX} Security*
516 *20)*, pages 1589–1604, 2020.
- 517 N. Singer. Microsoft urges congress to regulate use of facial recognition. *The New York Times*, 2018.
- 518 R. Singh, A. Agarwal, M. Singh, S. Nagpal, and M. Vatsa. On the robustness of face recognition algo-
519 rithms against attacks and bias. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
520 volume 34, pages 13583–13589, 2020.
- 521 M. Wang and W. Deng. Mitigating bias in face recognition using skewness-aware reinforcement
522 learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
523 pages 9322–9331, 2020.
- 524 T. Wang, J. Zhao, M. Yatskar, K.-W. Chang, and V. Ordonez. Balanced datasets are not enough:
525 Estimating and mitigating gender bias in deep image representations. In *Proceedings of the IEEE*
526 *International Conference on Computer Vision*, pages 5310–5319, 2019.
- 527 Z. Wang, K. Qinami, I. C. Karakozis, K. Genova, P. Nair, K. Hata, and O. Russakovsky. Towards
528 fairness in visual recognition: Effective strategies for bias mitigation, 2020.
- 529 K. Weise and N. Singer. Amazon pauses police use of its facial recognition software. *The New York*
530 *Times*, Jul. 10 2020. URL [https://www.nytimes.com/2020/06/10/technology/amazon-](https://www.nytimes.com/2020/06/10/technology/amazon-facial-recognition-backlash.html)
531 [facial-recognition-backlash.html](https://www.nytimes.com/2020/06/10/technology/amazon-facial-recognition-backlash.html).
- 532 M. J. Wilber, V. Shmatikov, and S. Belongie. Can we still avoid automatic face detection? In *2016*
533 *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–9. IEEE, 2016.
- 534 M. B. Zafar, I. Valera, M. Gomez Rodriguez, and K. P. Gummadi. Fairness beyond disparate treatment
535 & disparate impact. *Proceedings of the 26th International Conference on World Wide Web*,
536 Apr 2017a. doi: 10.1145/3038912.3052660. URL [http://dx.doi.org/10.1145/3038912.](http://dx.doi.org/10.1145/3038912.3052660)
537 [3052660](http://dx.doi.org/10.1145/3038912.3052660).
- 538 M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi. Fairness constraints: Mecha-
539 nisms for fair classification. In *Proceedings of the 20th International Conference on Artificial*
540 *Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA*, vol-
541 ume 54 of *Proceedings of Machine Learning Research*, pages 962–970. PMLR, 2017b. URL
542 <http://proceedings.mlr.press/v54/zafar17a.html>.
- 543 M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi. Fairness constraints: A flexible
544 approach for fair classification. *Journal of Machine Learning Research*, 20(75):1–42, 2019. URL
545 <http://jmlr.org/papers/v20/18-262.html>.
- 546 R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork. Learning fair representations. volume 28 of
547 *Proceedings of Machine Learning Research*, pages 325–333, Atlanta, Georgia, USA, 17–19 Jun
548 2013. PMLR. URL <http://proceedings.mlr.press/v28/zemel13.html>.
- 549 Z. Zhang, Y. Song, and H. Qi. Age progression/regression by conditional adversarial autoencoder. In
550 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5810–5818,
551 2017.

552 **Checklist**

- 553 1. For all authors...
- 554 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
555 contributions and scope? [Yes]
- 556 (b) Did you describe the limitations of your work? [Yes]
- 557 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 558 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
559 them? [Yes]
- 560 2. If you are including theoretical results...
- 561 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 562 (b) Did you include complete proofs of all theoretical results? [N/A]
- 563 3. If you ran experiments...
- 564 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
565 mental results (either in the supplemental material or as a URL)? [Yes]
- 566 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
567 were chosen)? [Yes]
- 568 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
569 ments multiple times)? [N/A]
- 570 (d) Did you include the total amount of compute and the type of resources used (e.g., type
571 of GPUs, internal cluster, or cloud provider)? [Yes]
- 572 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 573 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 3.
- 574 (b) Did you mention the license of the assets? [Yes] See Section 3.
- 575 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 576 (d) Did you discuss whether and how consent was obtained from people whose data you’re
577 using/curating? [Yes] See Section 3, and references within each of the papers that
578 introduce the datasets that we use and the noise models that we use.
- 579 (e) Did you discuss whether the data you are using/curating contains personally identifiable
580 information or offensive content? [Yes] See Section 3.
- 581 5. If you used crowdsourcing or conducted research with human subjects...
- 582 (a) Did you include the full text of instructions given to participants and screenshots, if
583 applicable? [N/A]
- 584 (b) Did you describe any potential participant risks, with links to Institutional Review
585 Board (IRB) approvals, if applicable? [N/A]
- 586 (c) Did you include the estimated hourly wage paid to participants and the total amount
587 spent on participant compensation? [N/A]