# Towards causal benchmarking of bias in face analysis algorithms

G. Balakrishnan[†‡]        Y. Xiong[‡]        W. Xia[‡]        P. Perona[*‡]

† Massachusetts Institute of Technology
∗ California Institute of Technology
‡ Amazon Web Services

## Abstract

Measuring algorithmic bias is crucial both to assess algorithmic fairness, and to guide the improvement of algorithms. Current methods to measure algorithmic bias in computer vision, which are based on *observational* datasets, are inadequate for this task because they conflate algorithmic bias with dataset bias.

To address this problem we develop an *experimental* method for measuring algorithmic bias of face analysis algorithms, which manipulates directly the attributes of interest, e.g., gender and skin tone, in order to reveal causal links between attribute variation and performance change. Our proposed method is based on generating synthetic "transects" of matched sample images that are designed to differ along specific attributes while leaving other attributes constant. A crucial aspect of our approach is relying on the perception of human observers, both to guide manipulations, and to measure algorithmic bias.

Besides allowing the measurement of algorithmic bias, synthetic transects have other advantages with respect to observational datasets: they sample attributes more evenly, allowing for more straightforward bias analysis on minority and intersectional groups, they enable prediction of bias in new scenarios, they greatly reduce ethical and legal challenges, and they are economical and fast to obtain, helping make bias testing affordable and widely available.

We validate our method by comparing it to a study that employs the traditional observational method for analyzing bias in gender classification algorithms. The two methods reach different conclusions. While the observational method reports gender and skin color biases, the experimental method reveals biases due to gender, hair length, age, and facial hair.
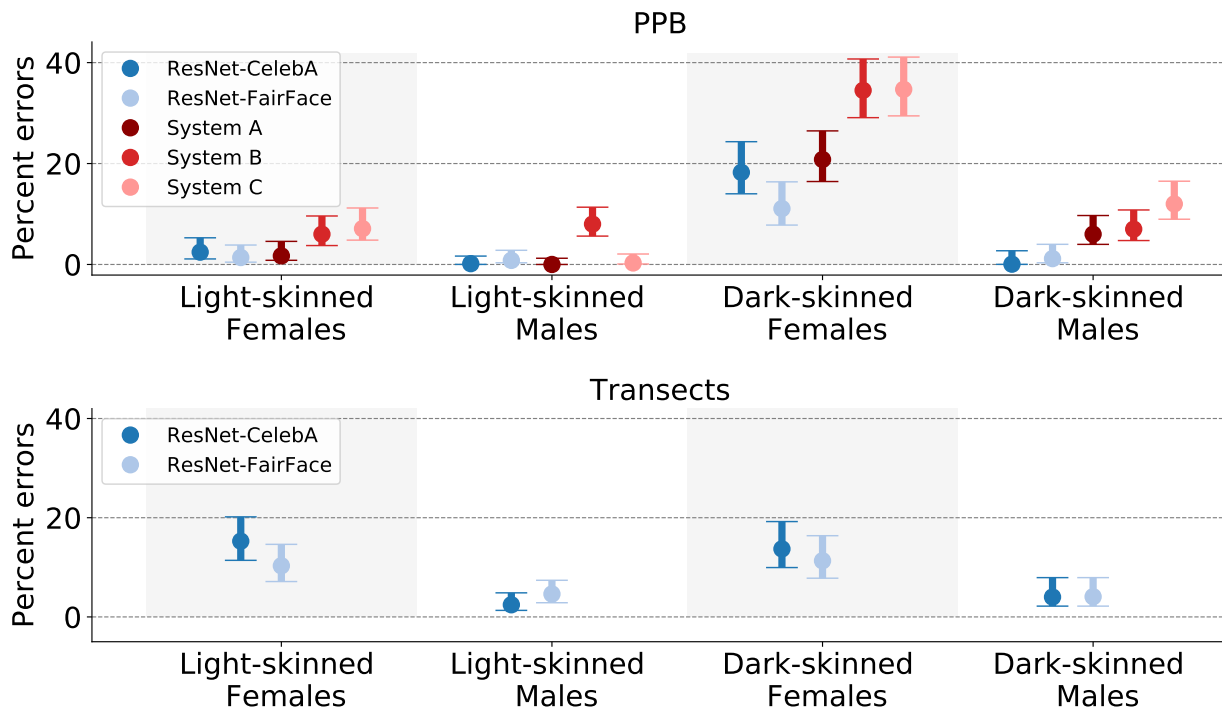
Figure 1: **Algorithmic bias measurements are test set dependent**. (Top) Gender classification error rates of three commercial face analysis systems (System A–C) were measured in 2017 on the Pilot Parliaments Benchmark (PPB) [12], an observational dataset of portrait pictures downloaded from the web sites of six national parliaments in Scandinavia and Africa. Error rates for dark-skinned females were found to be significantly higher than for other groups. We observed the same qualitative behavior when we replicated the study by training a standard classifier (ResNet-50) on two publicly available face datasets (CelebA, FairFace) and testing the two models thus obtained on a replica of the PPB dataset. (Bottom) Our experimental investigation using the Transects dataset, where sample faces are matched across attributes, reveals a different picture of algorithmic bias (see Fig 13, Sec. 5, and 6 for a more complete analysis).

# 1 Introduction

Automated systems trained using machine learning methods are increasingly used to support decisions in industry, medicine and government. While performance of such systems is often excellent, accuracy is not guaranteed, and needs to be assessed through careful measurements. Measuring *biases*, i.e., performance differences, across protected attributes such as age, sex, gender, and ethnicity, is particularly important for decisions that may affect peoples' lives. Unlike systems based on human judgment, where measuring and correcting biases is notoriously difficult, measuring and mitigating algorithmic bias is feasible and may become a powerful agent of progress towards more fair, accountable and transparent institutions [44, 55].

The prevailing technique for measuring the performance of algorithms is to measure statistics like error frequencies on a test set that is sampled *in the wild*, hopefully mirroring some of the data statistics that will be encountered in the field. Studies of algorithmic bias in computer vision [12, 10, 43, 47] have adapted this approach by adding one additional step: each image of the test set is

annotated for attributes of interest (e.g., ethnicity, gender and age), and the test set is then split into groups that have homogeneous attribute values. Comparing error rates across such groups yields predictions of bias. As an example, Fig. 1-top shows the results of a recent study of algorithmic bias in gender classification of face images. This type of study is called *observational*, because the independent variables (e.g., skin color and gender) are sampled from the environment, rather than controlled by the investigator.

Algorithmic bias is measured for two reasons. First, fairness: would changing a protected attribute, all else being equal, cause a systematic change in the output of the algorithm? For example, would two job applicants, that differed only by their gender or ethnicity, face predictably different outcomes [8]? The second reason for measuring bias is getting rid of it: which actions should one take to best improve the system's performance? For example, should the engineers who are in charge of developing systems A, B, and C (Fig. 1, top) infer that the best strategy is to add more examples of dark-skinned women to their training set? Thus, measuring algorithmic bias ultimately has one goal: revealing causal connections between attributes of interest and algorithmic performance.

Unfortunately, observational studies are ill-suited for drawing such conclusions. When one samples data in the wild, other variables may correlate with the variable of interest, and any one of the correlated variables may have an influence on the performance of the algorithm. Thus, it is difficult to impute the cause of performance differences to variations in the variable of interest – as the old saying goes: *"correlation does not imply causation."*

One simple instance of this problem is sample bias: samples in the wild may fail to represent specific combinations of variables of interest [38, 39, 40]. For example, the appearance of the parliamentarians in the PPB dataset [12] tends to be gender-stereotypical, e.g., very few males have long hair and almost no light-skinned females have short hair (Fig. 12, and [56]). The fact that hair length (a variable that may affect gender classification accuracy) is correlated in PPB with skin color (a variable of interest) complicates the analysis. In addition, the sample dataset that is used to measure bias is often not representative of the population of interest. For example, the middle-aged Scandinavians and Africans of PPB are not representative of, say, the broad U.S. Caucasian and African-American population [50]. While observational methods do yield useful information on disparate impact within a given test set population, generalizing observational performance predictions to different target populations is hit-or-miss [76] and can negatively impact underrepresented, or minority populations [53, 71]. In a nutshell, one would want a method that systematically identifies algorithmic bias while transcending the peculiarities of specific test sets.

Scientists in biology, medicine and the social sciences are well aware of this problem and have developed practices to discover, and to control for, confounding variables. A powerful approach to discovering cause-effect relationships is the *experimental method* which involves artificially manipulating the variable of interest, while fixing all the other inputs [8, 58]. This is not easy in the case of image data, leading us to ask the question: *Can one systematically measure bias in computer vision algorithms using the experimental method?* While this is not immediately intuitive [56], we find that the answer is yes, and offer a practical way forward.

Our approach (Fig. 2) generates the test images synthetically, rather than sampling them from the wild, so that they are varied selectively along attributes of interest. This is enabled by recent
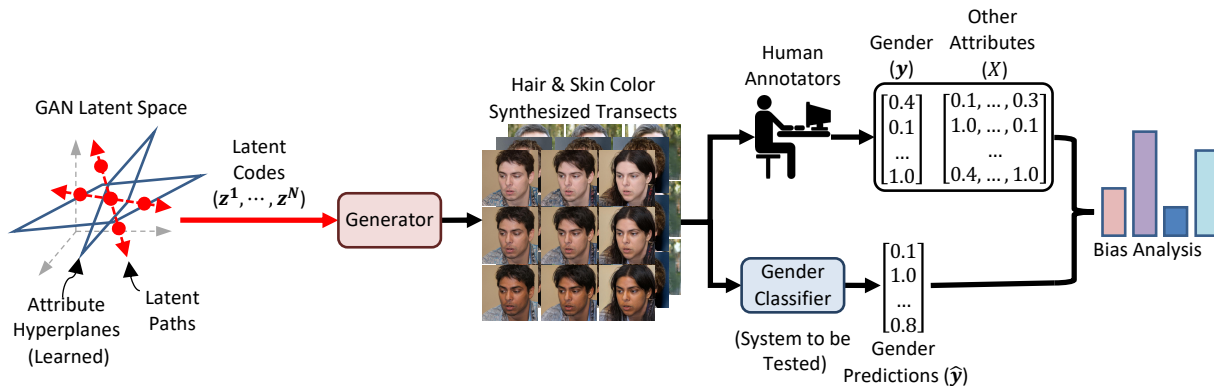
Figure 2: **Synopsis of our approach.** A generative adversarial network (Generator) is used to synthesize "transects," or grids of images, modifying selected attributes on synthetic faces (in this example: hair length and skin tone). This is accomplished by traversing the generator's latent space in attribute-specific directions. These directions are learned using randomly sampled faces and human annotators (not shown). Human annotations on the transects provide generator-independent ground truth to be compared with algorithm output to measure algorithm errors. Attribute-specific bias measurements are obtained by comparing the algorithm's predictions with human annotations as the attributes are varied. The depicted example may study the question: *Does hair length, skin tone, or any combination of the two have a causal effect on classifier errors?* Transects exploring other attributes are shown in Fig. 3, 4, and 9(a). The GUIs for human image annotation are shown in Fig. 6. Samples of image annotations are shown in Fig. 7.

progress in controlled and realistic image synthesis [36, 37], along with methods for collecting large amounts of accurate human annotations [11] to quantify the perceptual effect of image manipulations. Our synthesis approach can alter multiple attributes at a time to produce grid-like matched samples of images we call *transects*. We quantify the image manipulations with detailed human annotations which we then compare with algorithm output to estimate algorithmic bias.

We evaluate our methodology with experiments on two gender classification algorithms. We first find that our transect generation strategy creates significantly more balanced data across key attributes compared to "in the wild" face datasets. Next, inspired by [12], we use this synthetic data to explore the effects of various attributes like skin color, hair length, age and perceived gender on gender classifier errors. Our findings reveal that using an experimental method can change the picture of algorithmic bias (Fig. 13), which will affect the strategy of algorithm improvement, particularly concerning groups that are often underrepresented in training and test sets.

We view our work as a first step in developing experimental methods for algorithmic bias testing in computer vision which, we argue, are necessary to achieve trustworthy and actionable measurements. Much remains to be done, both in design and experimentation to achieve broadly-applicable and reliable techniques. In Sec. 6 we discuss limitations of the current method, and next steps in this research area.

## 2   Related Work

Benchmarking in computer vision has a long history [6, 9, 22] including face recognition [47, 60, 61, 62, 27, 28] and face analysis [12]. Some of these studies examine biases in performance, i.e., error rates across variation of important parameters (e.g. racial background in faces). Since these

studies are purely observational, they raise the question of whether the biases they measure depend on algorithmic bias, or on correlations in the test data. Our work addresses this question.

A dataset is said to be biased when combinations of features of interest are disproportionately represented or, equivalently, when such features are correlated. Computer vision datasets are often found to be biased [64, 76]. Human face datasets are particularly scrutinized [2, 20, 43, 45, 46, 54] because methods and models trained on these data can end up being biased along attributes that are protected by the law [44]. Approaches to mitigating dataset bias include collecting more thorough examples [54], using image synthesis to compensate for distribution gaps [46], and example resampling [48].

The machine learning community is active in analyzing biases of learning models, and how one may train models where bias is mitigated [3, 14, 18, 31, 33, 41, 46, 51, 68], usually by ensuring that performance is equal across certain subgroups of a dataset. Here we ask a complementary question: we assume that the system to be benchmarked is *pre-trained* and fixed, and we ask how to reliably measure algorithmic bias in pre-trained black-box algorithms.

Studies of face analysis systems [12, 43, 51] and face recognition systems [29, 47] attempt to measure bias across gender and skin-color (or ethnicity). However, the evaluations are based on observational rather than interventional techniques – and therefore any conclusions from these studies should be treated with caution. A notable exception is a recent study [56] using the experimental method to investigate the effect of skin color in gender classification. In that study, skin color is modified artificially in photographs of real faces to measure the effects of differences in skin color, all else being equal. However, the authors observe that generalizing the experimental method to other attributes, such as hair length, is too onerous if one is to modify existing photographs. Our goal is to develop a generally applicable and practical experimental method, where *any* attribute may be studied independently.

Recent work uses generative models to explore face classification system biases. One study explores how variations in pose and lighting affect classifier performance [2, 45, 46]. A second study uses a generative model to synthesize faces along particular attribute directions [19]. These studies rely on the strong assumption that their generative models can modify one attribute at a time. However, this assumption relies on having unbiased training data, which is almost always not practical. In contrast, our framework uses human annotations to account for residual correlations produced by our generative model.

Finally, there is research into interpreting neural networks. One strategy is to determine regions of the input that are salient, either through analysis of gradients or perturbations of the input image [13, 16, 24, 26, 69, 72, 74]. Network dissection approaches explore how particular neurons within a network affect the output, particularly in a semantic way [7, 82]. Testing with Concept Activation Vectors (TCAV) [42] provides explanations at a high level using directional derivatives to reveal the "conceptual sensitivity" of a model's prediction of a class (e.g., Smiling) to a concept. In contrast, our approach uses a synthesis model to create carefully modified input images, and human annotations to precisely quantify them.

# 3 Face Attribute Annotation in Synthetic Images

The face images used in our experiments are synthetic, and therefore there is no real person behind each image. Thus, there is no intrinsic ground truth for face attributes such as gender, hair length, and skin tone. Such attributes are instead established by human annotators. We clarify here what we mean when we talk about face attributes in the absence of a physical ground truth.

Many attributes have both intrinsic and extrinsic manifestations. For example, "emotion" may be studied at three levels [4]: an unconscious physiological state, conscious self-perception (feelings), and emotional display (e.g. facial expression) [17]. These quantities are *intrinsic* to a person's or an animal's body and are not directly accessible to a machine. By contrast, an *extrinsic* description, i.e., the report by an onlooker of his/her perception, are more easily accessible, and this is what the machine is trained to predict.

Since we are using synthetic images, it should be clear that we are not attempting to access the intrinsic state of a person: there is no person, and there is no intrinsic gender, ethnicity, age or emotion. However, perception of such attributes is possible. This is the same way that onlookers instinctively classify the *Venus of Milo* as "female" and Michelangelo's *David* as "male," despite the fact that they are idealized marble representations, rather than real people.

Thus, when we refer to the "age" or "gender" or any other attribute that is computed by a face analysis system from a picture, what we mean is *the algorithm's prediction of a casual observer's report of their perception of the outwards display of that attribute*. This is a bit of a mouthful, and that's why we use the abbreviated expression of "attribute," "age" or "gender." The attributes we measure from human observers are reports of subjective perceptions. However, as we find in Sec. 4.3, these measurements are consistent and reproducible across different observers, and so we consider statistics of such reports as objective quantities.

In our study, we discretize continuous face attributes. We have used six classes of age and skin tone, five of hair length, facial expression and gender, etc. (see Figs. 6 and 8). This choice was made to conform with the literature, e.g., the Fitzpatrick scale of skin tone [23], and to accommodate the abilities of non-expert casual observers, the "common person," whose perception we rely on in our experiments. We make no claim to have the perfect discretization scheme; other discretization choices may be better suited in different contexts.

Gender deserves a special mention: *gender identity* is often modeled as multi-dimensional [21]. However, here we are measuring *reports of gender perception* (an extrinsic variable), rather than gender identity (the intrinsic variable), and our subjects could not reliably report beyond the traditional one-dimensional M/F dimension. Therefore, following [12] we settled for one dimension, which we discretized into five steps to accommodate different levels of confidence and ambiguity.

# 4 Method

Our framework consists of two components: a technique to synthesize sets of images with control over semantic attributes, and a procedure using these synthesized images, along with human annotators, to perform analysis of a recognition system.

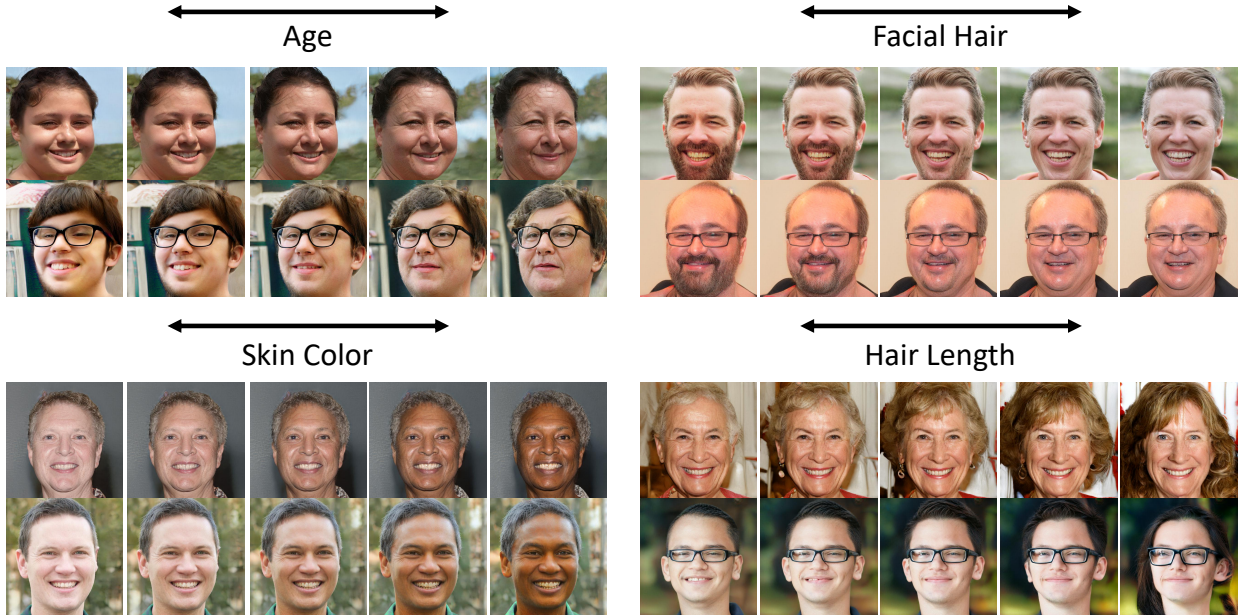In Sec. 4.1 we present our technique for attribute-controlled image synthesis. We introduce the

Figure 3: **1D transects.** $1 \times 5$ sample transects synthesized by our method for various attributes. Orthogonalization was used (see Fig. 5).

concept of *transect*, a grid-like construct of synthesized images with a different attribute manipulated along each axis. A transect gives control over the joint distribution of synthesized attributes allowing us to generate *matched samples* across multiple attributes, unlike related methods that operate on only one or two attributes at a time [19, 70, 73, 80]. We then collect human annotations for each transect image, to precisely quantify our modifications.

In Sec. 4.2 we present analyses we can perform using the annotated transects. We report a classifier's error rate, stratified along subgroups of a sensitive attribute. We also return a covariate-adjusted estimate of the *causal effect* of a binary attribute on the classifier's performance.

## 4.1 Transects: A Walk in Face Space

We assume a black-box face generator $G$ that can transform a latent vector $\mathbf{z} \in \mathcal{R}^D$ into an image $I = G(\mathbf{z})$, where $p(\mathbf{z})$ is a distribution we can sample from. In our study, $G$ is the generator of a pre-trained, publicly available state-of-the-art GAN ("StyleGAN2") [36, 37]. GAN latent spaces typically exhibit good disentanglement of semantic image attributes. In particular, empirical studies show that each image attribute often has a direction $\mathbf{v} \in \mathcal{R}^D$ that predominantly captures its variability [36, 82]. We base our approach on a recent study [82] for single attribute traversals in GAN latent spaces. That method trains a linear model to predict a particular image attribute from $\mathbf{z}$, and uses the model to traverse the $\mathbf{z}$-space in a discriminative direction. Our method generalizes this idea to synthesize image grids, i.e., *transects*, spanning arbitrarily many attributes.
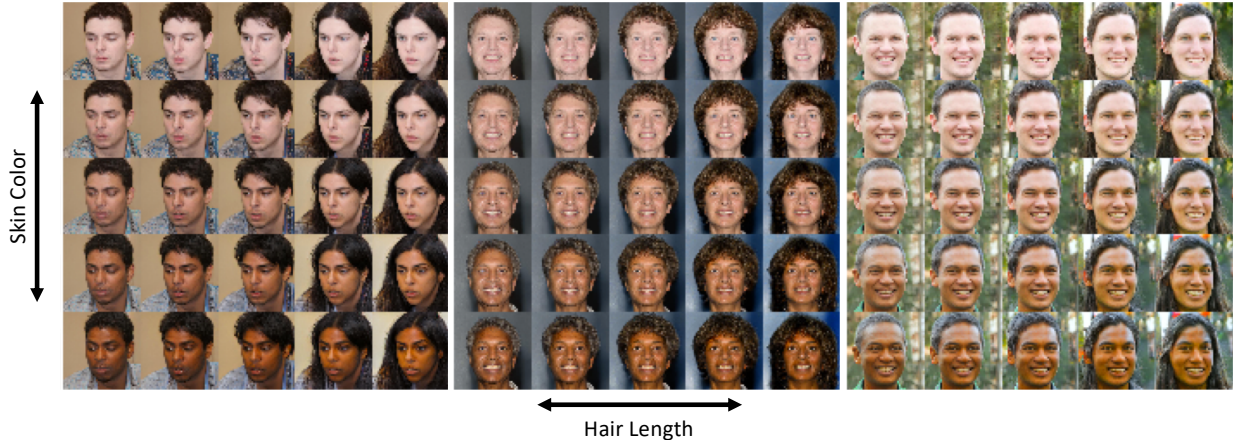
Figure 4: **2D transects.** $5 \times 5$ transects varying simultaneously hair length and skin tone. Multidimensional transects allow for intersectional analysis, i.e. analysis across the joint distribution of multiple attributes. Orthogonalization was used (see Fig. 5).

### 4.1.1 Estimating Latents-to-Attributes Linear Models

We first sample the latent space, measure the attributes at each location through human observers, and use these measurements to calculate principal axes of variation for attributes. More formally, let there be a list of $N_a$ image attributes of interest (age, gender, skin color, etc.). As explained below, we generate an annotated training dataset $\mathcal{D}_{\mathbf{z}} = \{\mathbf{z}^i, \mathbf{a}^i\}_{i=1}^{N_z}$, where $\mathbf{a}^i$ is a vector of scores, one for each attribute, for generated image $G(\mathbf{z}^i)$. The score for attribute $j$, $\mathbf{a}_j^i$, may be continuous in $[0, 1]$ or binary in $\{0, 1\}$.

We produce $\mathcal{D}_{\mathbf{z}}$ as follows. First, we sample a generous number of values of $\mathbf{z}^i$ from $p(\mathbf{z})$. Second, we obtain labels $\mathbf{a}^i$ from human annotators. A related study obtains labels by only processing the generated images through a trained classifier [82]. We generally avoid this approach because any biases of the classifier due to attribute correlations — precisely the phenomena we are trying to avoid — will leak into our method.

For each attribute $j$, we use $\mathcal{D}_{\mathbf{z}}$ to compute a $(D-1)$-dimensional linear hyperplane $h_j = (\mathbf{n}_j, b_j)$, where $\mathbf{n}_j$ is the normal vector and $b_j$ is the offset. For continuous attributes like age or skin color, we train a ridge regression model [34]. For binary attributes we train a support vector machine (SVM) classifier [15].

When sampling from StyleGAN2 using the native latent Gaussian distribution, we noticed a bias towards generating Caucasian-looking faces which is not surprising given the fact that it was trained on Flickr-Faces-HQ (FFHQ) – a public dataset that is skewed towards that demographic (see Fig. 11). However, using human annotations our method is able to partially mitigate this bias by directing sampling towards the relevant portions of the latent space (see following sections), so that it could generate a diversity of attributes. Nevertheless, training face synthesis GANs with a more diverse set of faces will be an important step in making our method more easily applicable.

### 4.1.2 Multi-attribute Transect Generation

The attribute hyperplanes may now be used to sample faces that vary along specific attributes. More formally: the hyperplane $h_j$ specifies the subspace of $\mathcal{R}^D$ with boundary or neutral values of attribute $j$, and the normal vector $\mathbf{n}_j$ specifies a direction along which that attribute primarily varies. To construct a one-dimensional, length-$L$ transect for attribute $j$, we first start with a random point $\mathbf{z}^i$ and project it onto $h_j$. We then query $L-1$ evenly-spaced points along $\mathbf{n}_j$, within fixed distance limits on both sides of the $h_j$. Fig. 3 presents some single transect examples (with orthogonalization, a concept introduced in the next section). We give further details on querying points in Sec. 4.1.4.

The 1D transect does not allow us to explore the joint space of several attributes, or to fix other attributes in precise ways when varying one attribute. We generalize to $K$-dimensional transects in Algorithm 1 to address this. The main extensions are: (1) we project $\mathbf{z}^i$ onto the intersection of $K$ attribute hyperplanes, and (2) we move in a $K$-dimensional grid in $\mathbf{z}$-space (see Fig. 4). Input $\mathbf{c}_k$ is a vector of decision values with respect to the hyperplane $h_k$, and $\mathbf{v}_k$ is a direction vector (equivalent to $\mathbf{n}_k$ here, until orthogonalization is introduced in the next section).

We are unaware of a simple closed-form solution to project $\mathbf{z}^i$ onto the intersection of arbitrarily many hyperplanes. We instead take an iterative approach: we sequentially project the point onto each hyperplane, and repeat this process for some number of iterations. Repeated projections onto convex sets, the hyperplanes in our case, is guaranteed to converge to a location on the intersection of the sets [81] which, in our case, is a single point. If the hyperplanes are perfectly orthogonal, this process converges in exactly one iteration; we empirically found convergence in fewer than 50 iterations.

### 4.1.3 Orthogonalization of Traversal Directions

The hyperplane normals $\{\mathbf{n}_j\}_{j=1}^{N_a}$ are not orthogonal to one another. If we set the direction vectors equal to these normal vectors in Algorithm 1, i.e., $\mathbf{v}_j = \mathbf{n}_j$, we will likely observe unwanted

---

**Algorithm 1:** $K$-attribute transect generation

**Input**: Generator $G$, tuples $\{(L_k, \mathbf{n}_k, b_k, \mathbf{v}_k, \mathbf{c}_k)\}_{k=1}^K$, where $L_k$ is a transect dimension, $(\mathbf{n}_k, b_k)$ is a hyperplane, $\mathbf{v}_k$ is a direction vector, and $\mathbf{c}_k$ are signed decision values.

**Output**: A $L_1 \times \cdots \times L_K$ transect $T^i$.

$\mathbf{z}^i \sim p(\mathbf{z})$
$\mathbf{z}^{i,0} = $ projection of $\mathbf{z}^i$ onto intersection of $\{(\mathbf{n}_k, b_k)\}_{k=1}^K$
**for** $l_1 = 1 \cdots L_1$ **do**
$\quad \vdots$
$\quad$ **for** $l_K = 1 \cdots L_K$ **do**
$\quad\quad T^i(l_1, \cdots, l_K) = G(\mathbf{z}^{i,0} + \sum_{k=1}^K \frac{\mathbf{c}_k[l_k]}{\langle \mathbf{v}_k, \mathbf{n}_k \rangle} \frac{\mathbf{v}_k}{\|\mathbf{v}_k\|})$

---

correlations between attributes. We reduce this effect by producing a set of modified direction vectors such that $\mathbf{v}_j \perp \mathbf{n}_k, \forall k \neq j$ (see Algorithm 2).

Fig. 5 illustrates the effects of orthogonalization for hair length and skin color. Without orthogonalization, the hair length transects exhibit unwanted changes in gender, with shorter hair also causing faces to appear more masculine. With orthogonalization, these unwanted changes are removed. In contrast, we see no clear difference in skin color transects with and without orthogonalization, indicating that the skin color hyperplane was already near-orthogonal to the other attribute hyperplanes.

### 4.1.4   Setting Step Sizes and Transect Dimensions

If human annotation cost were negligible, we could simply query many grid locations with large transect dimensions $L$ to capture subtle appearance changes over the dynamic ranges of the attributes. But given constrained resources, we set $L$ to small values. For example, $L = 5$ for the 1D transects in Fig. 3 and 2D transects in Fig. 4, and $L = 2$ for the 3D transects in Fig. 9. For each attribute $j$, we manually set min/max signed decision values with respect to $h_j$, and linearly interpolate $L_j$ points between these extremes to obtain $\mathbf{c}_j$. We set per-attribute min/max values so that transects depict a full dynamic range for most random samples.

## 4.2   Analyses Using Transects

We assume a target attribute of interest, e.g., gender, and a target attribute classifier $C$. We will use transect images to perform bias analysis on $C$. Though an ideal transect will modify only selected attributes at a time, in practice, unintended attributes may also be accidentally modified. In addition, the degree to which an attribute is altered varies across transects. To measure and control for these factors we annotate each image of each transect, resulting in a second dataset $\mathcal{D}_{transect} = \{I^i, \mathbf{a}^i\}_{i=1}^{N_{images}}$ of images and human annotations.

We denote the ground truth gender of image $I^i$ (as reported by humans) by $y^i$, and $C$'s prediction by $\hat{y}^i$. For ease of analysis, we discretize the remaining attributes into bins, and assign

---

**Algorithm 2:** Orthogonalization

**Input**: Vectors $\{\mathbf{n}_j\}_{j=1}^{N_a}$.
**Output**: Vectors $\{\tilde{\mathbf{n}}_j\}_{j=1}^{N_a}$, where $\tilde{\mathbf{n}}_j \perp \mathbf{n}_k, \forall k \neq j$

$Q, R \leftarrow$ QR-factorization of matrix $[\mathbf{n}_1, \mathbf{n}_2, \cdots, \mathbf{n}_{N_a}]$
**for** $i = 1 \cdots N_a$ **do**
$\quad$ $\tilde{\mathbf{n}}_i = \mathbf{n}_i$
$\quad$ **for** $j = 1 \cdots N_a$ **do**
$\quad\quad$ **if** $i \neq j$ **then**
$\quad\quad\quad$ $\tilde{\mathbf{n}}_i = \tilde{\mathbf{n}}_i - \frac{Q_j \cdot \langle Q_j, \tilde{\mathbf{n}}_i \rangle}{\langle Q_j, Q_j \rangle}$

---

Figure 5: **1D transects with and without orthogonalization.** Without orthogonalization (Sec. 4.1.2), decreasing hair length results in more masculine-looking faces. This phenomenon is not as apparent after orthogonalization (Sec. 4.1.3). We see only slight orthogonalization differences in the skin color transects.

an independent binary variable to each bin [25]. For instance, we may represent the 'skin color' attribute with six binary variables, corresponding to the six levels shown in Fig. 6 (top right). For a given image, only one of these six variables would be set to $1$ – often called a 'one-hot encoding.' We denote the vector of concatenated binary covariates for image $i$ by $\mathbf{x}^i$, and the classification error by $e^i = \ell(\hat{y}^i, y^i)$, where $\ell(\cdot, \cdot)$ is an error function.

Our first analysis strategy is to simply compare $C$'s error rate across different subgroups in the population. Let $E_j^s$ denote the average error of $C$ over test samples for which covariate $j$ is equal to $s \in \{0, 1\}$:

$$E_j^s = \frac{\sum_i e^i \mathbb{1}(\mathbf{x}_j^i = s)}{\sum_i \mathbb{1}(\mathbf{x}_j^i = s)}. \tag{1}$$

If the data is generated from a perfectly randomized or controlled study, the quantity $E_j^1 - E_j^0$ is a good estimate of the "average treatment effect" (ATE) [5, 32, 57, 66] of covariate $j$ on $e$, or the average change in $e$ over all examples when covariate $j$ is flipped from $0$ to $1$, with other covariates fixed. For example, the ATE of the "dark skin" covariate captures the average change in $C$'s error when each person's skin tone is changed from non-dark to dark. Exactly computing the ATE from an observational dataset is not possible, because we do not observe the counterfactual case(s) for each data point, e.g., the same person with both light and dark skin tones.

Though our transects come closer to achieving an ideal controlled study than do observational datasets "from the wild" (see Sec. 5.3 for empirical validation), there may still be some confounding between covariates in practice (see Fig. 18 for an example). Since any observable confounder may be annotated in $\mathcal{D}_{transect}$, we can employ covariate-adjusted ATE estimators [63, 65, 78]. One simple covariate adjustment approach is to train a linear regression model predicting $e^i$ from $\mathbf{x}^i$:

$$e^i = \epsilon^i + \beta_0 + \sum_j \beta_j \mathbf{x}_j^i, \tag{2}$$

where $\beta$'s are parameters, and $\epsilon^i$ is a per-example noise term. $\beta_j$ captures the ATE, the average change in $e$ given one unit change in $\mathbf{x}_j$ holding all other variables constant, provided: (1) a linear model is a reasonable fit for the relationship between the dependent and independent variables, (2) all relevant attributes are included in the model (i.e., no hidden confounders), and (3) no attributes that are influenced by $\mathbf{x}_j$ are included in the model, otherwise these other factors can "explain away" the impact of $\mathbf{x}_j$.

An experimenter can never be completely sure that (s)he has satisfied these conditions but (s)he can strive to do so through careful consideration. Discretizing and binarizing attributes helps with (1), though we still found some non-linear influences between covariates in our experiments (see Sec. 5.5.1). As an example of (2), we found that earrings may be an important attribute that we did not account for in our analysis (see Fig. 19).

Finally, when the outcome lies in a fixed range, as is the case in our experiments with $e^i \in [0, 1]$, we use logistic instead of linear regression. $\beta_j$ then represents the expected change in the *log odds* of $e$ for a unit change in $\mathbf{x}_j$. We use such a logistic regression analysis in our experiments (see Sec. 5.5).

## 4.3   Human Annotation

We collect human annotations on the synthetic faces to construct $\mathcal{D}_{\mathbf{z}}$ and $\mathcal{D}_{transect}$. The annotators were recruited on Amazon Mechanical Turk [11] through the AWS SageMaker Ground Truth service [1]. Annotators evaluated each image for seven attributes: gender, facial hair, skin color, age, makeup, smiling, hair length and image fakeness. Each attribute was evaluated on a discrete scale. Each annotator evaluated each image for one attribute at a time. For each image, we collected 5 annotations per attribute for a total of 40 annotations per image.

    We discretized each attribute using three to six levels. For example, we use the Fitzpatrick six-point scale for skin color [23], and split age into six groups ranging from children to senior citizens. For complete details about subgroups for each attribute, along with samples of our Mechanical Turk survey layouts please see Fig. 6.

    The number of annotations that are needed by our method is rather formidable. However, we found that this is not an obstacle in practice. In our experiments, $\mathcal{D}_{\mathbf{z}}$ consists of 5,000 images, and $\mathcal{D}_{transect}$ consists of 1,000 8-image transects (see examples in Fig 9). The total number of annotations was thus 13,000 (images) x 8 (attributes) x 5 (annotations per image and per attribute) = 0.52M annotations. Amazon Mechanical Turk delivered on average 10-20 annotations per second, thus annotations took about 10 hours to complete over two separate sessions. Annotators were paid 1.2c per annotation, earning 10-15 US$ per hour.

    One may be concerned that annotators may not be able to give meaningful attribute annotations on synthetic images. Therefore we explored the level of agreement of annotator responses, both in a number of pilot experiments, and in the annotations we collected for the main experiment. Fig. 7 shows the raw annotations for one 1D transect and three attributes. One may see that there are very few outlier annotations, and that in most cases annotations fall in one or two neighboring attribute levels. Fig. 8 (top left) shows a distribution of per-image annotation standard deviations, split by attribute. One unit corresponds to the dynamic range of each attribute. For most attributes, the median annotator standard deviation is near 0.1, i.e. less than the separation between attribute levels. These observations indicate good agreement between annotators and suggest that annotations are meaningful and reproducible.

    Fig. 8 (top right) presents the distribution of mean annotator fakeness scores for the synthesized images. Only a small portion of images are deemed "Likely fake" or "Fake for sure." Realism of images is particularly important in our analysis, since image artifacts can unknowingly affect the decisions of gender classifiers. In our experiments, we remove images with a fakeness score above a certain threshold (see Sec. 5.2.1). Fig. 8 (bottom) shows example synthesized images organized by mean human fakeness score.

# 5   Experiments

In order to test our method on a practical application, we experiment with benchmarking bias of gender classifiers. The Pilot Parliaments Benchmark (PPB) [12], a dataset of faces of parliament members of various nations, was the first wild-collected test dataset to balance gender and skin color with the goal of fostering the study of gender classfication bias across both attributes. The
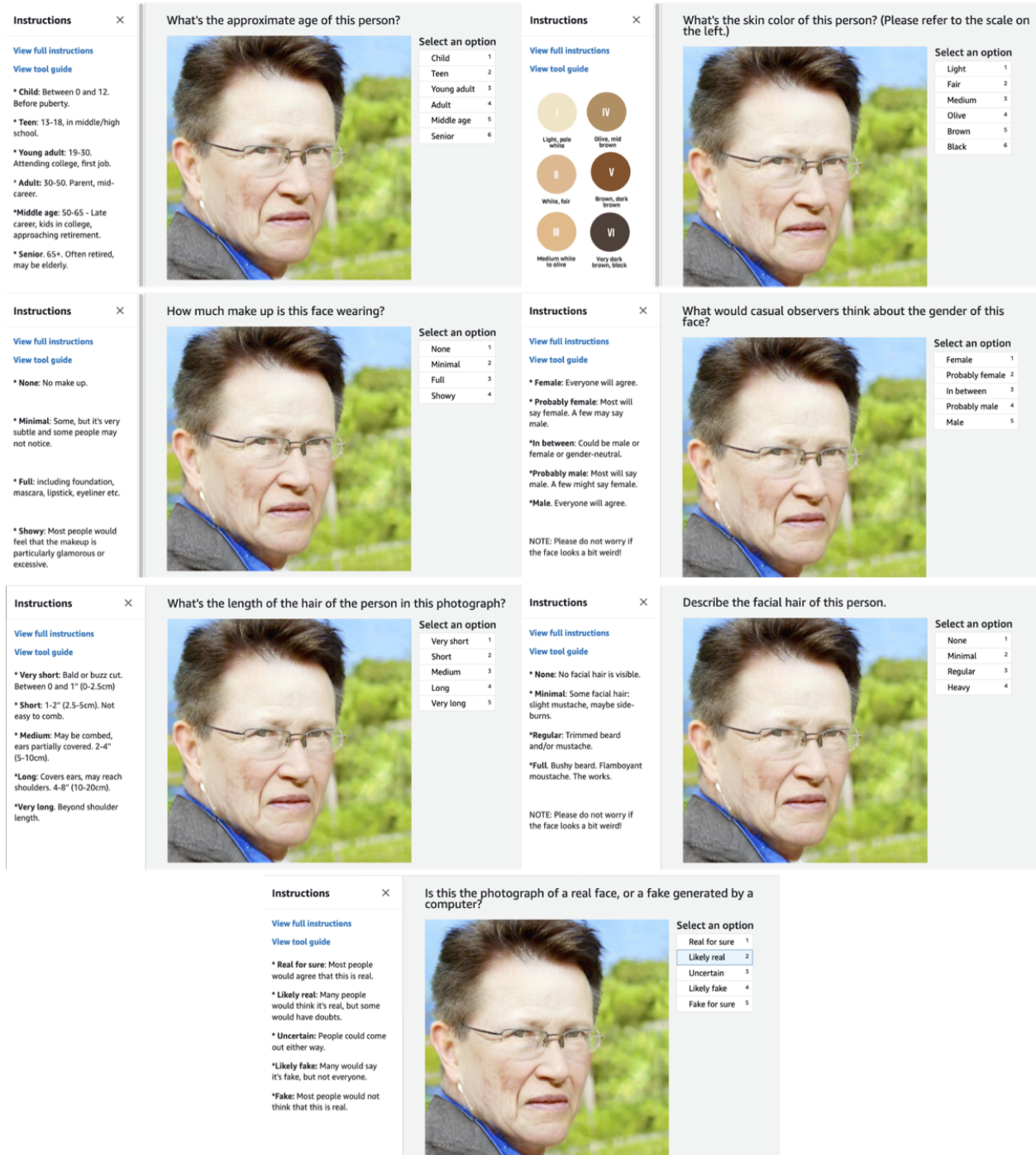
Figure 6: **Screenshots of the graphical user interface** for seven annotations we collected from Amazon Mechanical Turk annotators using the SageMaker Ground Truth service [1].

authors of that study found a much larger error rate on dark-skinned females, as compared to other groups and conjectured that this is due to bias in the algorithms, i.e., that the performance of the
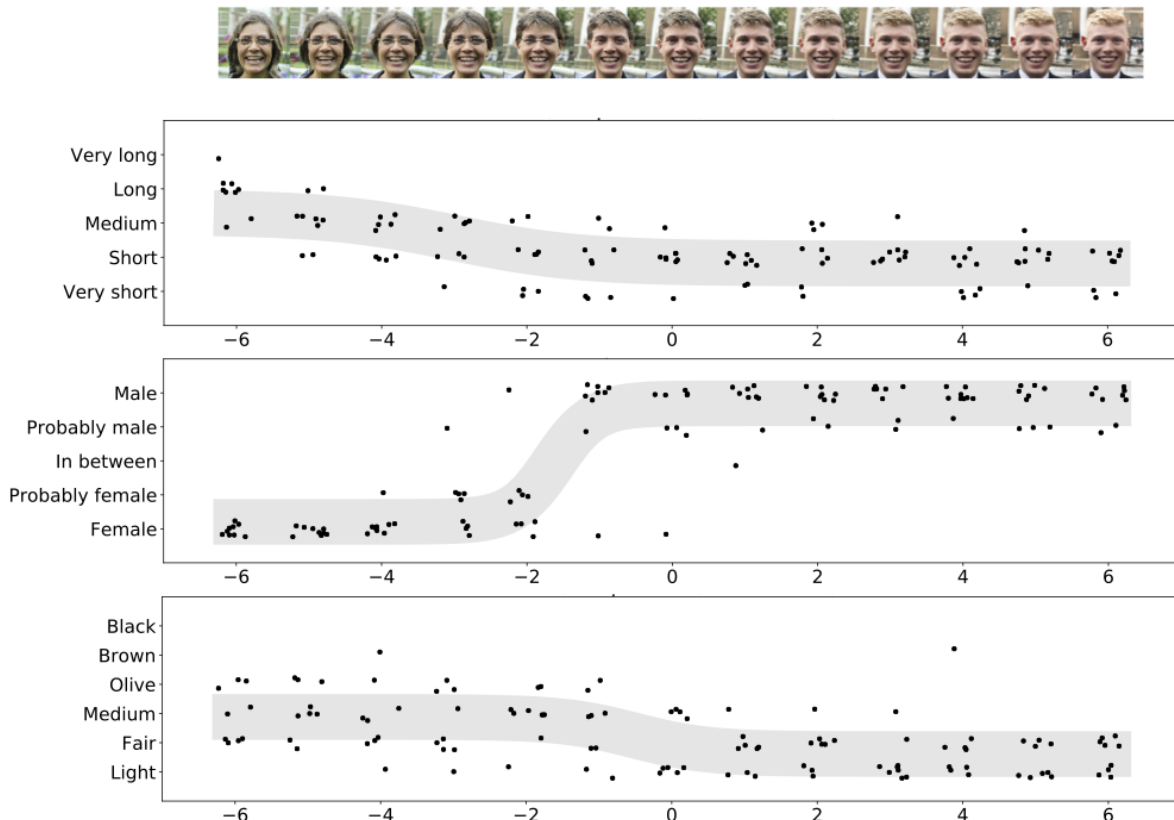
Figure 7: **Annotation consistency.** Hair length (top), gender (middle) and skin tone (bottom) annotations on a 13-image 1D transect. This transect was annotated in a pilot experiment to fine tune our GUIs and to evaluate the consistency of the annotators, and not used in our main experiment. Here nine annotations were obtained for each attribute and for each image. The annotations are shown as dots below each image. The $x$ axis increments one unit from one image to the next. A small amount of noise was added in $x$ and $y$ in order to visualize the individual annotations. The thick gray curves show the fit of a logistic function to the data. Annotations typically fall within one or two neighboring attribute levels. There are very few outliers. For a quantitative overall analysis see Fig.8.

algorithm changes when gender and skin color are changed, all else being equal. Our method allows us to test this hypothesis.

## 5.1 Gender Classifiers

We trained two research-grade gender classifier models, each using the ResNet-50 architecture [30]. The first was trained on the CelebA dataset [49], and the second on the FairFace dataset [35]. CelebA is the most popular public face dataset due to its size and rich attribute annotations, but is known to have severe imbalances [19]. The FairFace dataset was introduced to mitigate some of these biases.

We trained our classifiers for 20 epochs with the binary cross-entropy loss. We set the learning

**Fakeness**

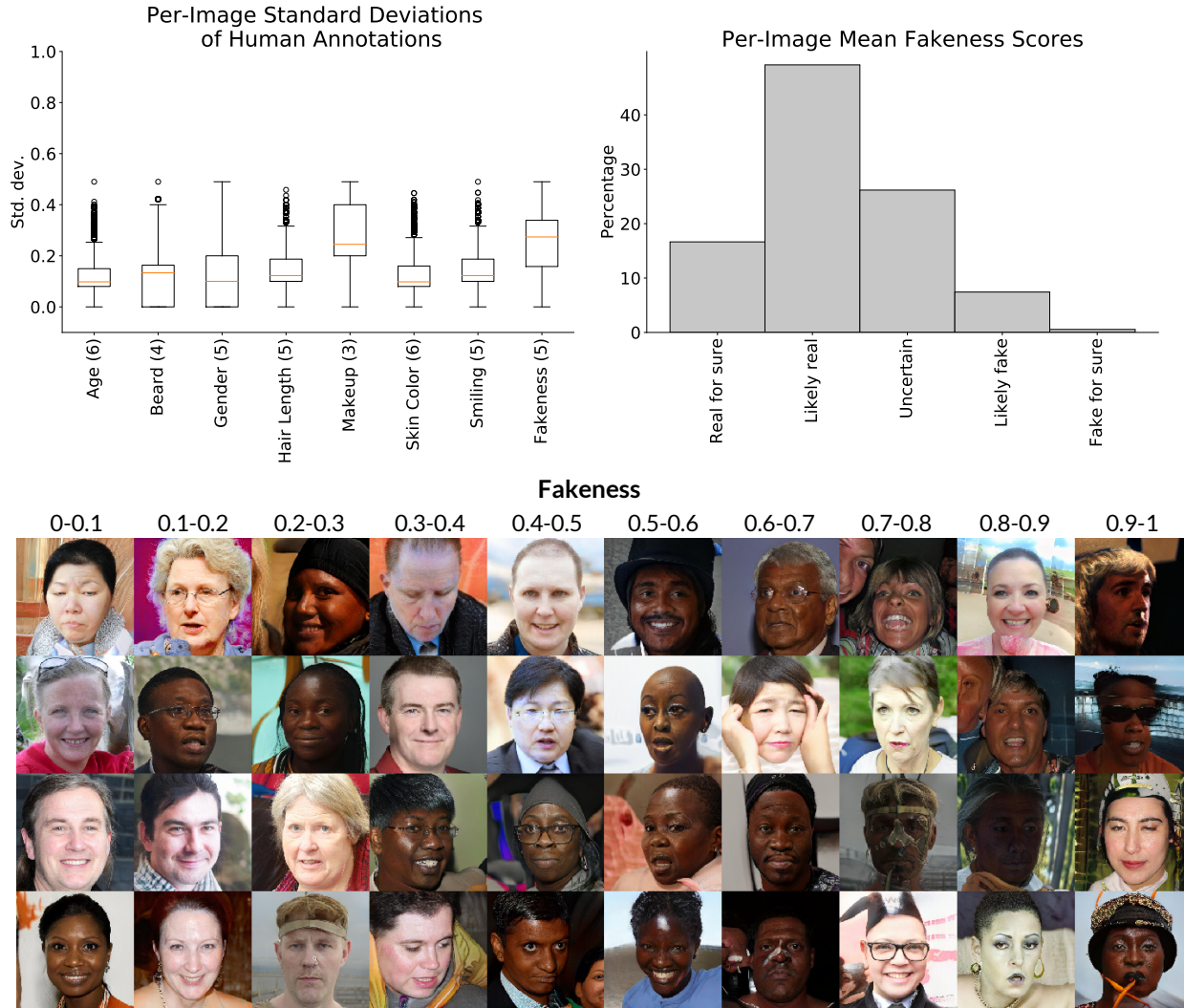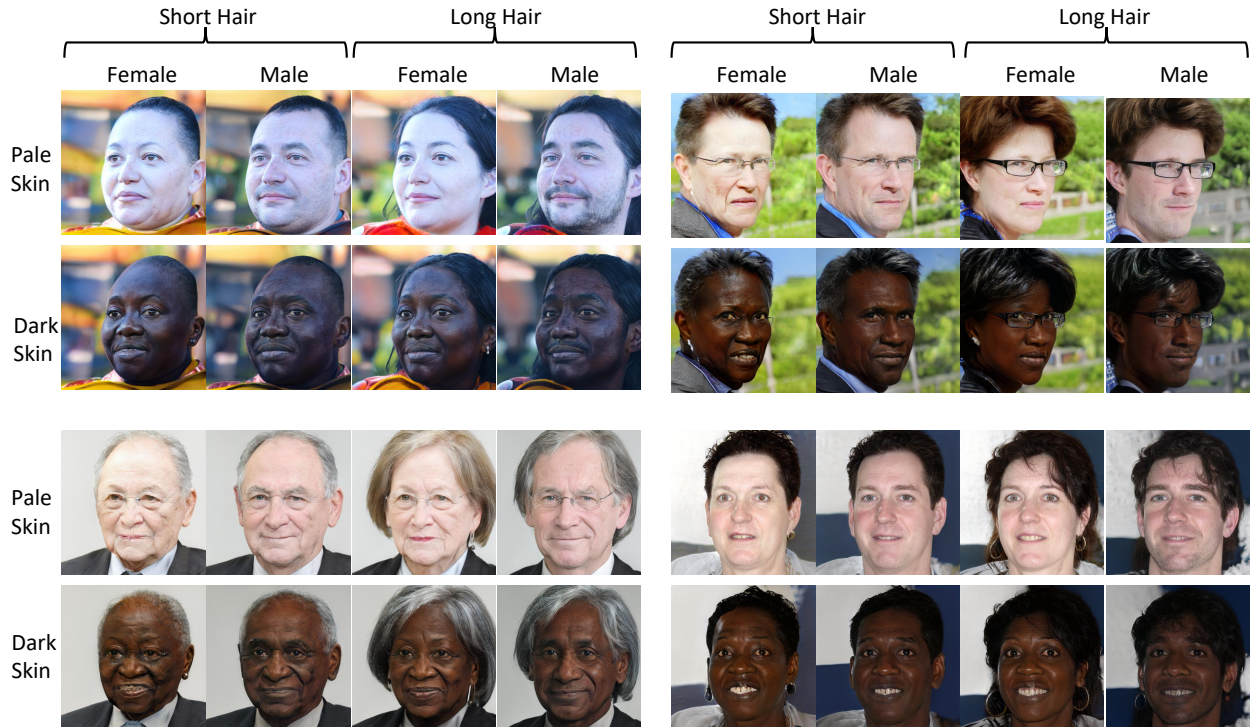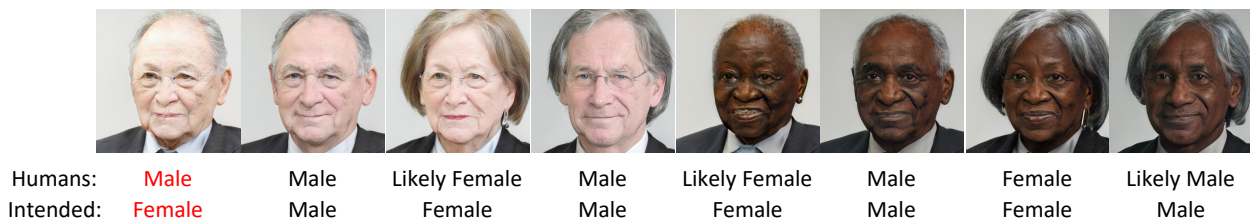| 0-0.1 | 0.1-0.2 | 0.2-0.3 | 0.3-0.4 | 0.4-0.5 | 0.5-0.6 | 0.6-0.7 | 0.7-0.8 | 0.8-0.9 | 0.9-1 |



Figure 8: **Annotation quality and image realism.** (Top left) Distributions of per-image standard deviations of human annotations for each of the attributes we considered (one unit = dynamic range of the attribute). Five annotators were asked to provide a rating for each attribute of each image. The number of rating options per attribute is indicated in brackets next to the attribute's name. The median standard deviations (red lines) are comparable to the quantization step, indicating good annotator agreement. (Top right) We asked our annotators to rate the realism of the images. The distribution of such scores is shown. Fewer than 10% of the ratings indicated fake or likely fake, suggesting that the synthetic images we randomly sampled are fairly realistic. (Bottom) we show examples of synthesized faces organized by mean human fakeness scores. Images with high fakeness scores were removed from the experiments (see Sec.5.2.1).

rate at $1e^{-4}$ for the first 10 epochs, and $1e^{-5}$ for the final 10 epochs. To avoid a baseline bias of predicting one gender over another, we enforced the likelihood of sampling male and female faces during training to be equal.

We decided not to test commercial system for two reasons. First, reproducibility — the models we test may be re-implemented and re-trained by other researchers at any time, while commercial systems are black boxes which may change unpredictably over time. Second, our ResNet-50

Short Hair | Long Hair | Short Hair | Long Hair
Female | Male | Female | Male | Female | Male | Female | Male

Pale Skin

Dark Skin

Pale Skin

Dark Skin

(a) **Examples of transects used in our experiments.**



Humans: Male | Male | Likely Female | Male | Likely Female | Male | Female | Likely Male
Intended: Female | Male | Female | Male | Female | Male | Female | Male

(b) **Human perception of the generators' manipulations.**

Figure 9: **Sample of 3-attribute transects used in our experiments.** We created 1,000 $2 \times 2 \times 2$ transects spanning skin color, hair length and gender – four examples are shown in (a). We set step sizes in such a way that we obtained pale-to-dark skin tones, short-to-medium hair lengths and M/F gender (see Sec. 4.1.4). Besides the intentionally modified attributes, other face attributes are held constant. For each image in each transect we collected human annotations to measure the perceived attributes. In (b) we show human-annotated gender values of the bottom-left transect in (a) side-by side with the generator's intended values. Humans label the first face as a male, though the generator intended to produce a female. In all our experiments we used human perception, rather than intended generator attributes, as the ground truth.

models show biases comparable to those observed in the original study by [12] (see Fig. 2.

17

## 5.2 Transect Data

To produce the synthetic images for our transects, we used the generator from the StyleGAN2 architecture trained on Flickr-Faces-HQ (FFHQ) [36, 37]. This generator has both a multivariate Normal input noise space, $\mathcal{N}(\mathbf{0}, \mathbf{I})$, as well as an intermediate "style space." To train the latent space linear models (see Sec. 4.1.1), we sampled $5000$ vectors from the noise distribution, and labeled the generated images with human annotators (see Sec. 4.3). However, we use the *style space* as the latent space in our method, because we found it better suited for disentangling semantic attributes. We trained linear regression models to predict age, gender, skin color and hair length attributes from style vectors. For the remaining attributes — facial hair, makeup and smiling — we found that binarizing the ranges and training a linear SVM classifier works best.

We generated 3D transects across subgroups of skin color, hair length, and gender following the procedure described in Sec. 4.1.2. We use a transect size of $2 \times 2 \times 2$, with grid decision values (specified by input vector $\mathbf{c}$ in Algorithm 1) spaced to generate a pale-to-dark transition along the skin color axis, short-to-medium length along the hair length axis, and male-to-female along the gender axis. We set the decision values by trial-and-error, and made them equal for all transects: $(-1.5, 1.7)$ for skin color, $(-0.5, 0)$ for hair length, and $(-1.75, 1.75)$ for gender. We generated $1000$ such transects, resulting in $8000$ total images. Fig. 9 presents four example transects. The general characteristics of the faces — besides the intentionally modified attributes — are held constant.

### 5.2.1 Dataset Pruning

Not all synthesized images are ideal for our analysis. Some elicit ambiguous human responses (Fig. 8 top left) or are unrealistic (Fig. 8 top right). Furthermore, others may not belong clearly to one of our two intended categories for the gender, hair length, and skin color attributes. We addressed these points by first removing any image with a mean fakeness score greater than or equal to "Likely fake" ($0.75$ in the normalized range of $[0, 1]$). We also removed faces with attribute values in the normalized subranges of $[0.4, 0.6]$ for skin color and gender, and $[0.3, 0.5]$ for hair length (see Fig. 10 for examples). After these pruning steps, we were left with 5713 images.

## 5.3 Comparison of Transects to Real Face Datasets

Fig. 11 analyzes attribute distributions for the CelebA-HQ, FFHQ and PPB datasets, along with our transects, stratified by gender. The wild-collected datasets contain significant imbalances across gender, particularly with hair length. They also have biases in age, with a larger percentage of males being older than females. An interesting correlation is that males are also more likely to smile in these data. In contrast, our transects exhibit more balance across gender. They depict more males with medium-to-long hair, and fewer females with very long hair. Our transects also have a bimodal skin color distribution, and an older population by design, since we are interested in mimicking those population characteristics of PPB. All datasets are imbalanced along the "Beard" and "Makeup" attributes — this is reasonable since we expect these to have strong correlations with gender.
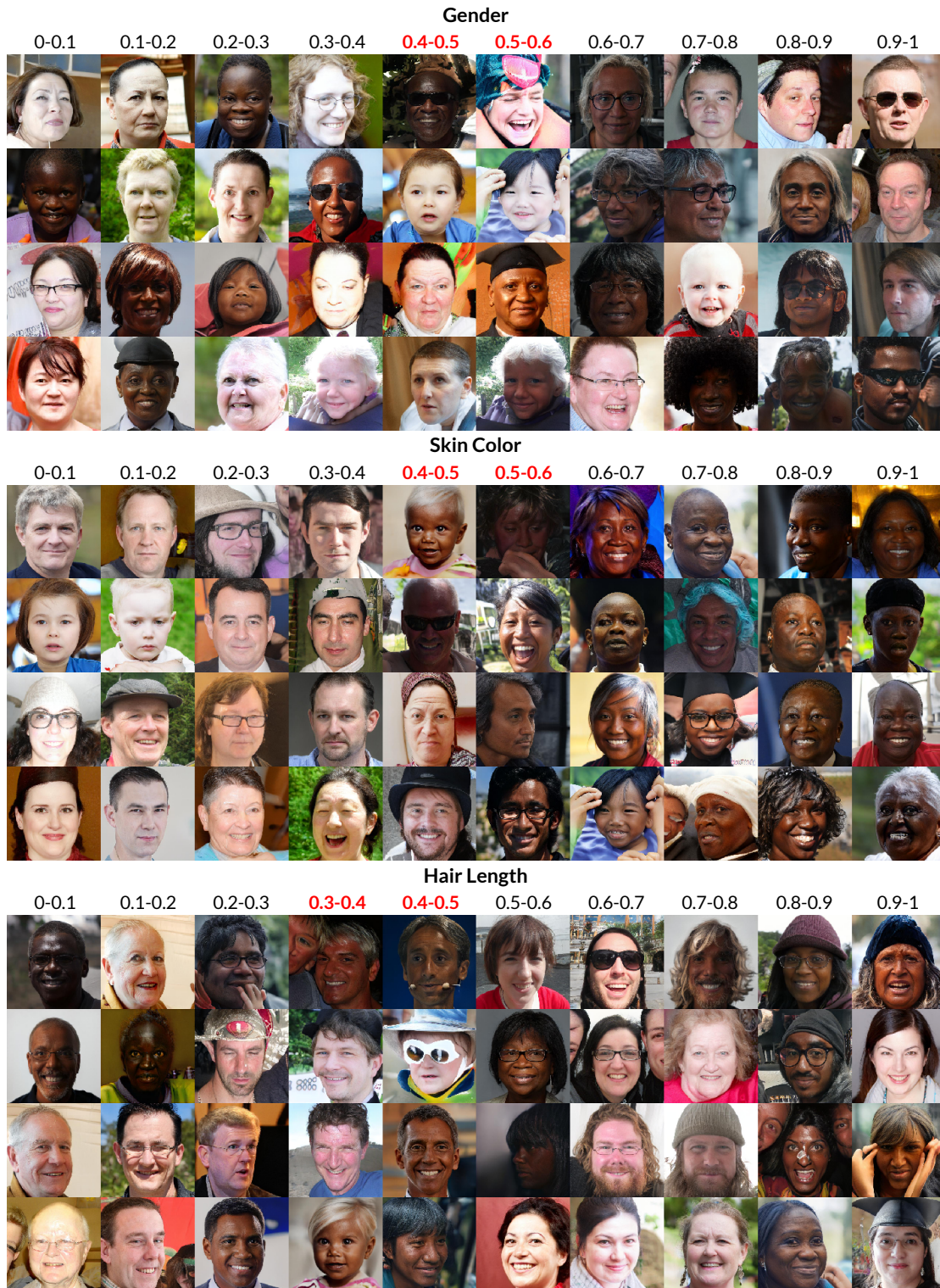
18

**Gender**

| 0-0.1 | 0.1-0.2 | 0.2-0.3 | 0.3-0.4 | <span style="color:red">0.4-0.5</span> | <span style="color:red">0.5-0.6</span> | 0.6-0.7 | 0.7-0.8 | 0.8-0.9 | 0.9-1 |

**Skin Color**

| 0-0.1 | 0.1-0.2 | 0.2-0.3 | 0.3-0.4 | <span style="color:red">0.4-0.5</span> | <span style="color:red">0.5-0.6</span> | 0.6-0.7 | 0.7-0.8 | 0.8-0.9 | 0.9-1 |

**Hair Length**

| 0-0.1 | 0.1-0.2 | 0.2-0.3 | <span style="color:red">0.3-0.4</span> | <span style="color:red">0.4-0.5</span> | 0.5-0.6 | 0.6-0.7 | 0.7-0.8 | 0.8-0.9 | 0.9-1 |

Figure 10: **Samples of synthesized faces, organized by mean human annotation scores.** In our analysis, we omitted faces from ranges indicated in red to focus on clearly perceived females/males, light/dark skin tones, and short/long hair lengths.
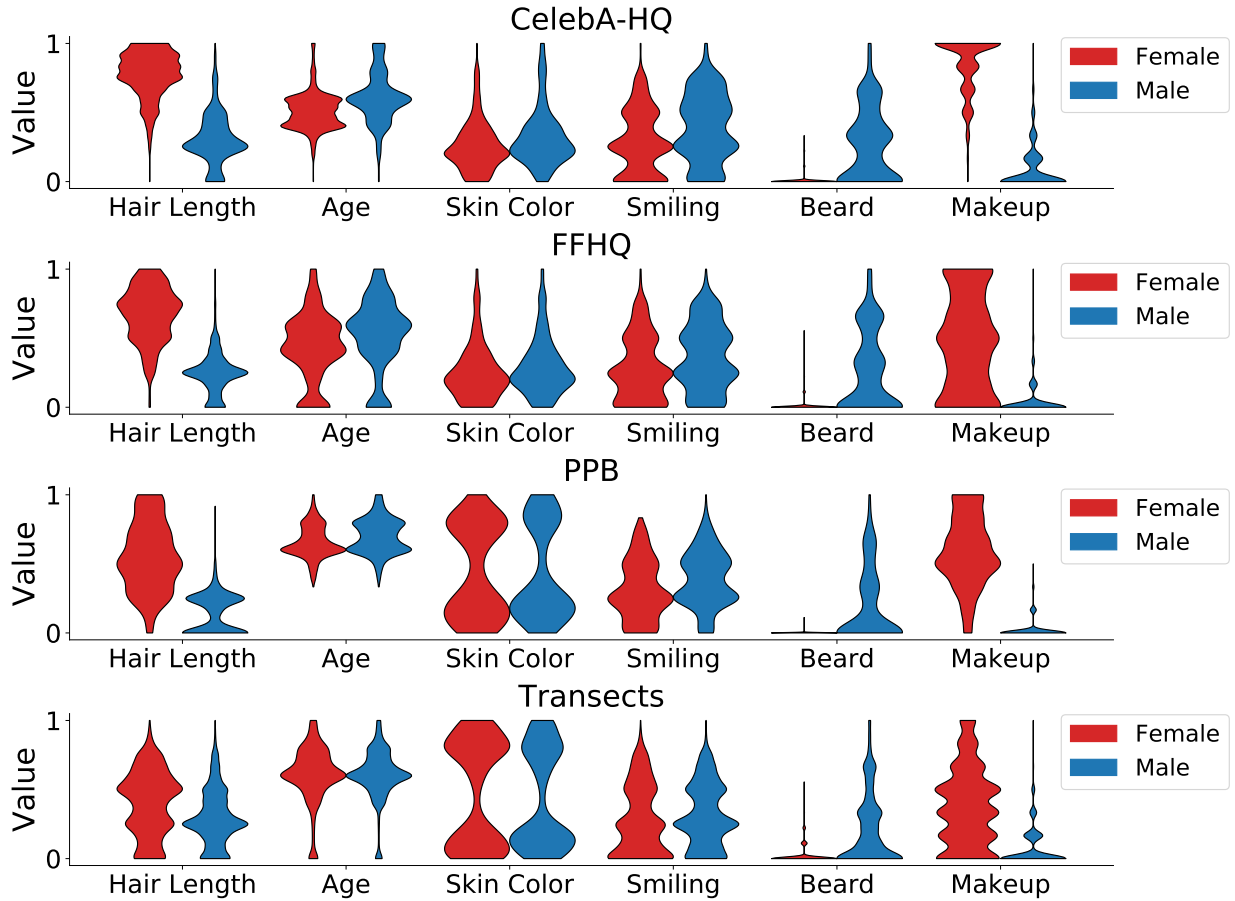
19

Figure 11: **Attribute distributions by dataset and gender groups.** "Violin" plot widths are proportional to frequency counts, and each violin is scaled so that its maximum count spans the full width. Wild-collected datasets have greater attribute imbalances across gender than synthetic transects, e.g. longer hair and younger ages for women. We designed our transects to mirror PPB skin color distribution and age distributions, while mitigating hair length imbalance. Hair length vs. skin color distributions are further explored in Fig.12.



Figure 12: **Hair length distributions by gender and skin color groups.** In the wild-collected datasets hair length is correlated with skin color, when gender is held constant. Synthetic transects may be designed to minimize this correlation.

In an ideal matched study, sets of images stratified by a sensitive attribute will exhibit the same distribution over remaining attributes. Put simply, no other attribute should be strongly correlated

with the attribute being manipulated. Fig. 12 stratifies by skin color. We see correlations of hair length distributions and skin colors in all the wild-collected data, while the synthetic transects exhibit much better balance.

## 5.4 Analysis of Bias

We now analyze the performance of the classifiers on PPB and our transects. We verify that the classifiers exhibit similar error patterns to the commercial classifiers already evaluated on PPB [12]. Because PPB only consists of adults, we remove children and teenagers (age $< 0.4$ in the normalized $[0, 1]$ scale) from our transects to make a more direct comparison, leaving us with 5335 total images.

Fig. 1 presents classification errors split by gender (M/F) and skin color (L/D). We replicated the reported errors of the commercial classifiers in [12], and report the errors of our classifiers on our in-house version of PPB. All classifiers perform significantly worse on dark-skinned females. Fig. 13 and Fig. 14 present classification errors, stratified by gender/hair length/skin color combinations. We can make a number of broad-stroke, qualitative observations:

- The broad pattern of errors is similar across PPB and transects, with more errors on the left (females) than on the right (males).

- Transect errors are either comparable or higher than in PPB, indicating that synthetic faces can be at least as challenging as real faces. Most significantly, errors are nonzero on males, which allows the study of relative difficulties when attributes are varied.

- In PPB, there are few males with long hair and few females with short hair and light skin, making measurements unreliable for these categories. This is not a problem with transects, where faces are matched by attributes.

- Transect errors are higher when hair is shorter for women. However, hair length has a negligible effect for males (see Fig. 18 for a possible explanation).

- There is no consistent transect error pattern in skin tone: within homogeneous groups changing skin tone does not seem to affect the performance of either algorithm. For example, females with long hair see no significant difference in classification error between light vs. dark skin. Looking at PPB alone, we could not make this observation, since skin tone is so strongly correlated with hair length.

## 5.5 Regression Analysis

In order to obtain a quantitative assessment of effect (or lack thereof) of attributes on classifier error, we investigated further by calculating covariate-adjusted causal effects. For each gender classifier model we trained an $L2$-regularized logistic regression model to predict that classifier's error conditioned on all attributes.
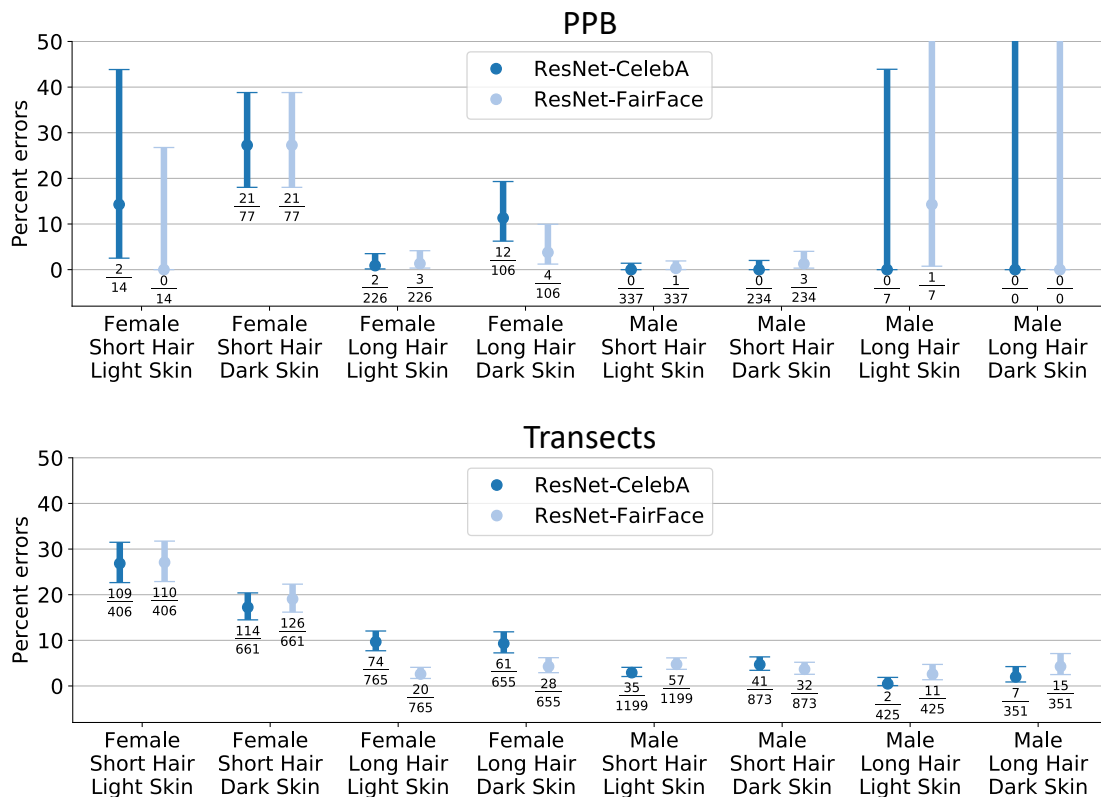
Figure 13: **Algorithmic errors, disaggregated by intersectional groups for wild-collected (PPB, top) and synthetic (transects, bottom).** Wilson score 95% confidence intervals [79] are indicated by vertical bars, and the misclassification count and total number of samples are written below each bar. PPB has few samples for several groups, such as short-haired, light-skinned females and long-haired males (see Fig. 12). Synthetic transects provide numerous test samples for all groups. The role of the different attributes in causing the errors is studied in Fig. 15.
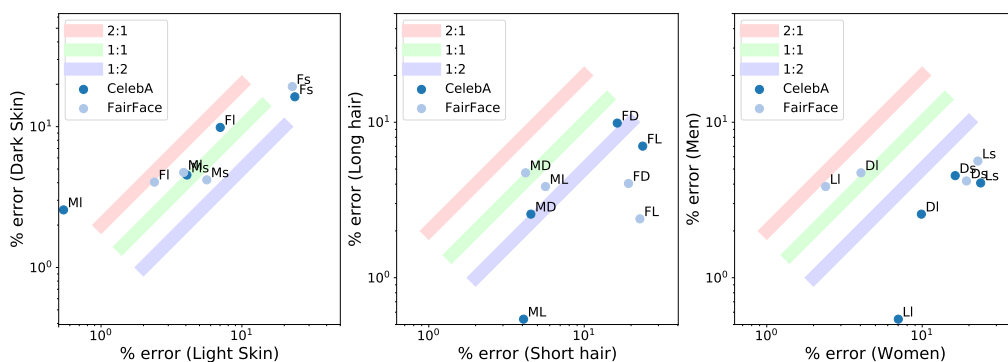


Figure 14: **Scatter plots of error rates using data from Fig 13 (transects).** Each dot compares the error rates of a pair of groups that differ by one attribute only (indicated in the label of the $x$ and $y$ axes). The two letters near each dot indicate the shared attributes ('M/F' indicate male and female, 'D/L' indicate dark and light skin, and 's/l' indicate short and long hair). Dots falling along the equal error line indicate that skin tone has little or no effect on error. In contrast, females and persons with short hair have higher error rates.
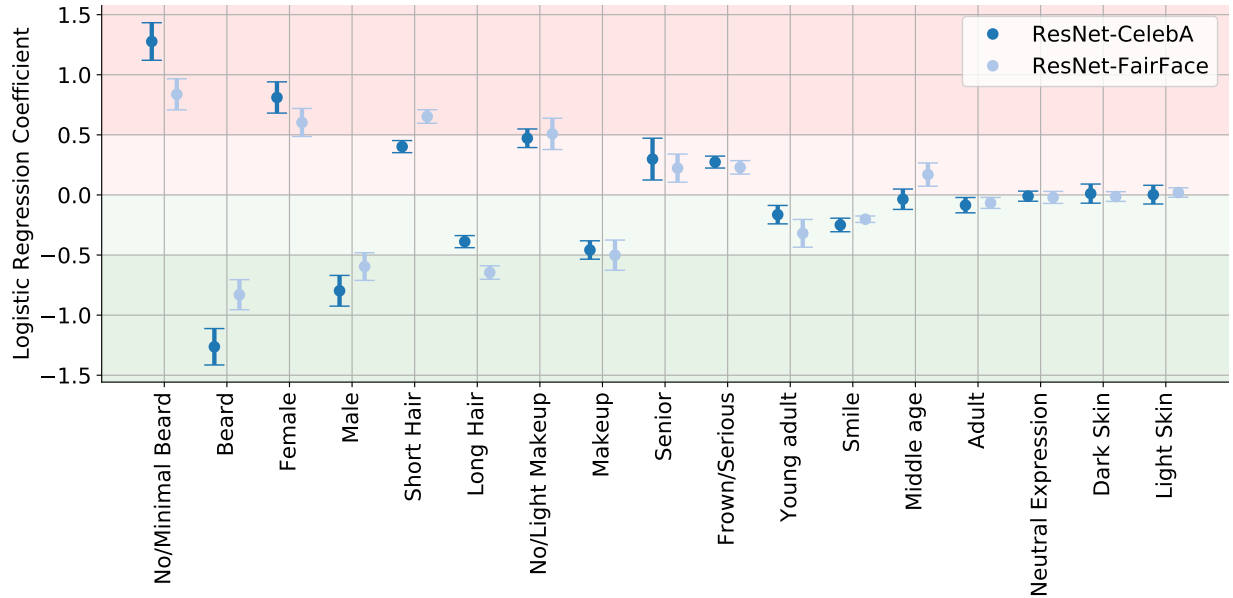
Figure 15: **Logistic regression coefficient values.** The logistic regression model is trained to predict *absolute errors* of the gender classifiers on our transect data given attributes as input. Coefficients represent the change in *log odds* of the error for a change of 1 unit of each attribute. Larger coefficient magnitudes indicate more important variables, and positive(red)/negative(green) values correspond to variables that increase/decrease classifier error. Each attribute subgroup labeled on the $x$-axis is represented by a binary variable in the regression model, and we order attributes in this plot from large-to-small coefficient magnitudes. Error bars report standard deviations that were computed via bootstrapping 1000 times.

We discretized attributes into levels, and assigned a binary variable to each level. We used the same discretization for hair length (short vs. long hair), skin color (light vs. dark skin) and gender (female vs. male) used in our experiments thus far. We used two levels for beard (no/light beard vs. beard) and makeup (no/light makeup vs makeup), three for facial expression (serious/frown vs. neutral vs. smile), and the original semantic levels for age described in Fig. 6. In all, this resulted in 17 input variables to our logistic regression model. We used scikit-learn's LogisticRegression function [59], and set the regularization parameter to 1.

Fig. 15 presents coefficients for both logistic regression models. Recall that each coefficient represents the change in log odds of the classifier's error for a change of 1 unit of each covariate (see Sec. 4.2). Error bars depict standard deviations, obtained by bootstrapping the dataset 1000 times. A person's facial hair, gender, makeup, hair length and age all have significant effects on classification error, and skin color has a negligible effect. Our main experimental conclusion is that observational (PPB) and experimental (transects) methods are fundamentally at odds on the causes of algorithm bias in gender classification algorithms. Observational analysis on wild-collected PPB suggests that a combination of gender and skin tone are implicated, while our experimental method using synthetic transects suggests that other attributes are far more important than skin tone.
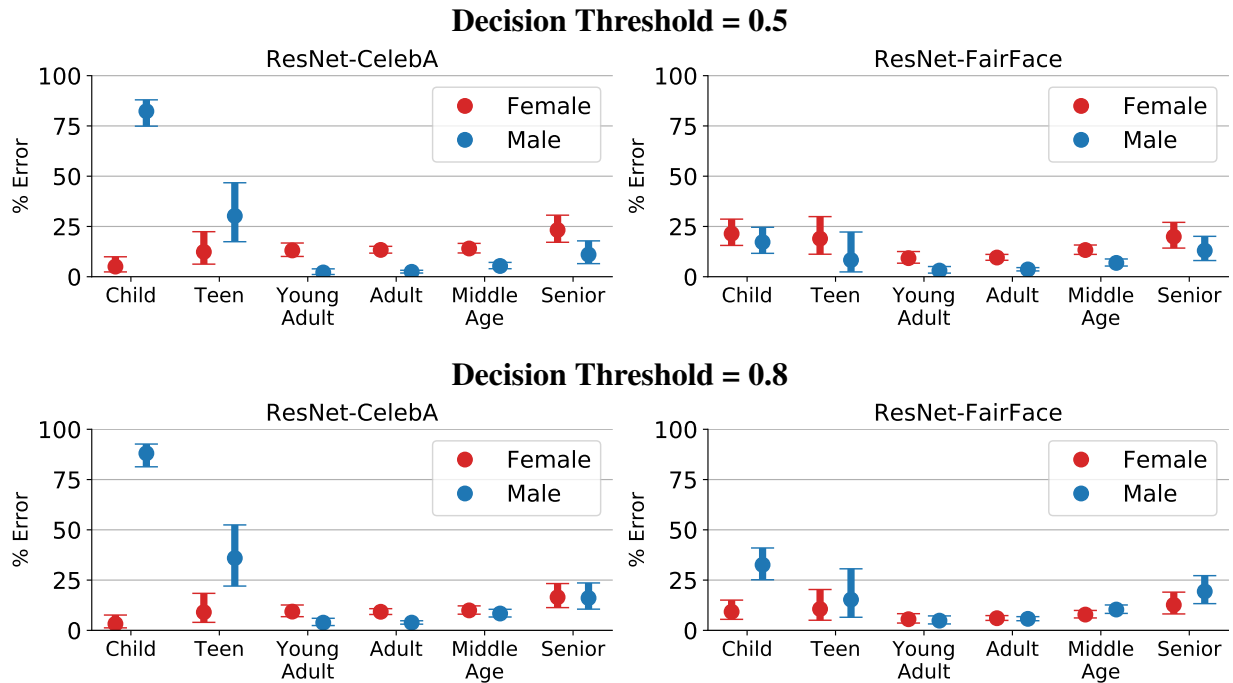
Figure 16: **Errors by gender and age group on our transect images.** The two top plots were obtained by using a decision threshold equal to 0.5, and show a prevalence of female errors. The bottom two plots were obtained with a threshold equal to 0.8, chosen to minimize overall error. There is a non-uniform influence of age on errors. Both models tend to have lower errors for young to middle-aged adults. The differences in errors between genders are fairly consistent for adults, but differ for children, teenagers and seniors, illustrating a combined age-gender bias in the algorithms.

### 5.5.1 Joint Effects of Attributes on Classification Error

Our regression analysis makes a simplifying assumption that each covariate has an independent, linear effect on classification error. The independence assumption can be a poor one. For example, Fig. 14-right shows that error rates vary across different intersectional groups of skin color and hair length in a way that is not simply a linear combination of each attribute.

This is also the reason we removed children and teenagers from our analysis, as these individuals tend to have different appearance characteristics from adults. Fig. 16 illustrates this, by breaking down error rates by age and gender subpopulations for two classifier decision thresholds. The difference in error rates between the genders is fairly consistent for young adults to middle-aged individuals, but vary for children/teenagers and seniors. This demonstrates that age and gender have joint effects on errors.

Fig. 17 shows faces from our synthesized transects on which the ResNet models were most incorrect. For each gender misclassification direction, we show faces on which the model predictions were farthest from the average human annotator response. ResNet-CelebA tends to heavily misclassify young male children/babies as female, in line with the quantitative result in Fig. 16.
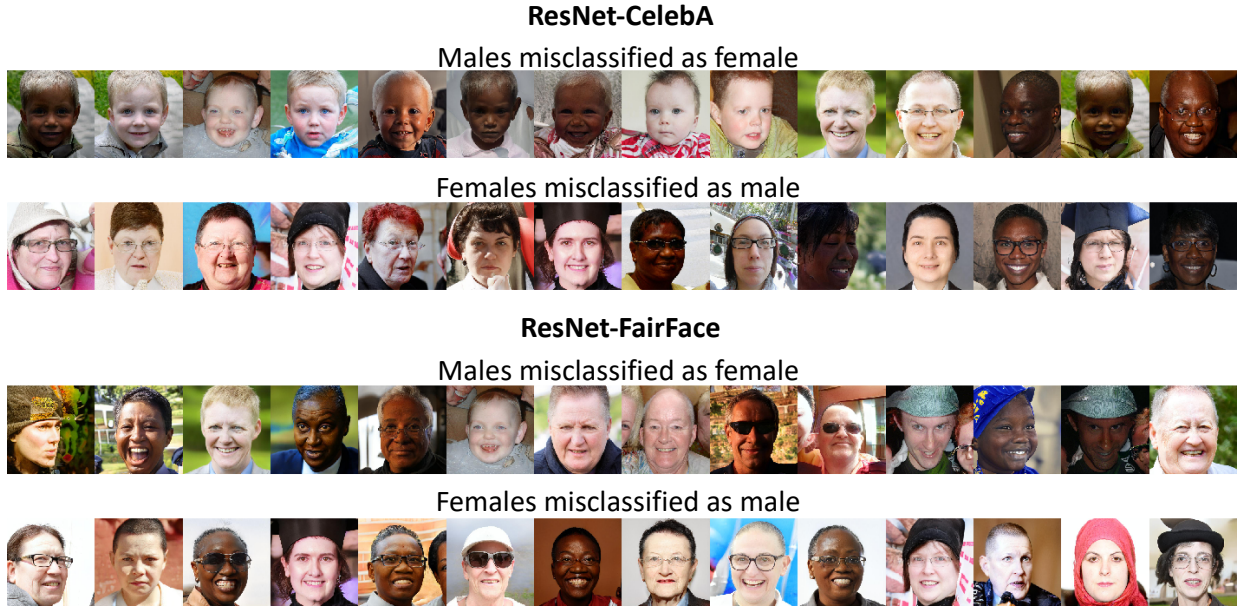
**ResNet-CelebA**

Males misclassified as female



Females misclassified as male



**ResNet-FairFace**

Males misclassified as female



Females misclassified as male



Figure 17: **Images with largest errors.** Synthetic faces on which the classifiers most deviated from the mean human annotations.

# 6 Discussion and Conclusions

## 6.1 Summary

Our study leads us to three main conclusions. First, the experimental approach to measuring algorithmic bias in computer vision is feasible. Second, the experimental approach may yield quite different conclusions from traditional observational studies. Third, when analyzing algorithmic bias, a broad spectrum of attributes and attribute combinations should be considered besides the ones of immediate interest. We examine each in detail below.

Our experimental approach is made possible by combining recent progress in image synthesis with detailed human annotations collected by crowdsourcing. Image synthesis, calibrated by human annotations, allows us to generate transects of matched samples, i.e., groups of images which differ only along attributes of interest. In contrast to previous attribute-specific methods [56] *any* attribute may be explored, provided that it can be annotated by humans. By relying on human ground truth annotations, one does not need to rely on the synthesis method being perfect.

The experimental method and our synthesis-based experimental approach, offer a number of attractive properties and advantages over traditional observational methods:

1. **Causal inferences on bias are possible.** Our method generates approximately matched samples across selected attributes, allowing for counterfactual analysis, e.g., *"Would the algorithm have made the same mistake if the same person had had a different skin color?"* Observational image data are almost never matched.

2. **Bias may be measured for underrepresented groups.** Image synthesis allows, to a great

Example 1                                      Example 2



Figure 18: **Correlated attribute modifications.** We found that our method sometimes adds a beard to a male face when attempting to only modify hair length. This is an example of an imprecise intervention which can complicate downstream bias analyses. This bias may be due to the training data itself (men with long hair tend to have facial hair), or injected by the algorithm.

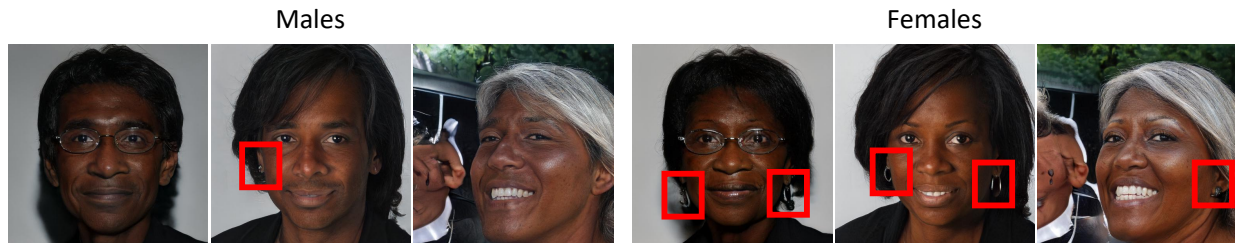Males                                          Females



Figure 19: **Hidden confounders.** There is always the possibility of a hidden confounder lurking in a dataset. As an example, we found — after already collecting annotations — that our method tends to add earrings when transitioning from dark-skinned men to dark-skinned women, a cue that a gender classifier might use to perform disproportionately well on the latter group. Because we did not annotate this attribute, it is hidden to our analysis. Interestingly, one male in this image also has an earring; that earring becomes larger for his female counterpart.

degree, uniform sampling of the space of attributes of interest — gender, skin color and hair length in our experiments. This is very difficult to do when one relies on images that are sampled from natural distributions, which tend to be long-tailed and therefore where some groups are underrepresented.

3. **Bias may be measured for intersectional groups.** Our method allows researchers to draw causal inferences across groups that are defined by specific attribute combinations. Single-attribute analysis may conceal biases affecting groups defined by the combination of multiple attributes [40]. Some such combinations are often vastly undersampled in natural data.

4. **Bias measurements are valid across different populations.** This is because the experimental method identifies causally linked attributes, independent of the prevalence of these attributes, i.e., the bias measurements are a property of the algorithm and not of the population on which it is used. By contrast, observational measurements do not generalize beyond the narrowly defined population where the data was collected. Furthermore, by combining appropriately the contribution of different attributes, one may predict the effects both of *disparate treatment* [52] and *disparate impact* [67] on a specific population.

5. **Accurate bias measurements may be made quickly and inexpensively.** Image synthesis is fast and inexpensive, and crowdsourced image annotation is also relatively fast and affordable. By contrast, assembling large datasets of natural images is laborious and expensive – it may take years and substantial investment, which may only be afforded by large organizations. Thus, synthetic data has the potential to democratize testing for bias.

6. **Ethical and legal concerns are greatly reduced.** Collecting face image datasets in the wild requires great care to respect the privacy and dignity of individuals, the rights of minors and other vulnerable groups, as well as copyright laws. By contrast, synthetic datasets are free from such risks because they do not depict real people.

The experimental analysis (transects) and traditional observational analysis (using PPB) diverged most significantly on the effect of skin color, which the observational study flagged as significant and the experimental method found to be not significant in determining algorithmic bias. The experimental method reveals a number of additional sources of bias: age, hair length and facial hair (Fig. 15). The two methods agree on gender. Our analysis suggests that the difference between the conclusions of the two methods is likely due to the correlation of hair length, skin color and gender in PPB (see Fig. 11 and Fig. 12). Consequently, if one does not control for hair length, the classifiers' bias towards assigning gender on the basis of hair length is read as a bias concerning dark-skinned women. The triple correlation between hair length, gender and skin color had been noticed in a previous study [56].

The main reason for measuring algorithmic bias is to get rid of it. Error and bias measurements guide scientists and engineers towards effective corrective measures for improving the performance of their algorithms. It is instructive to view the different predictions of the two methods through this lens. The correlational study based on PPB (Fig. 1) may suggest that, in order to reduce biases in our classifiers, more images of dark-skinned women should be added to their training sets. The experimental method leads engineers in a different direction. First, more training images of long-haired men and short-haired women of all races are needed. Second, correcting age bias requires more training images in the child-teen and, possibly, senior age groups.

Finally, a lesson from our study is that it is important to consider a rich number of attributes and attribute combinations, besides the one(s) of immediate interest. This is for two reasons. First, unobserved confounders can have strong effects and need to be included in the analysis. Second, the combined effect of attributes can be strongly nonlinear (see the interaction of age and gender in Fig. 16), and therefore an intersectional analysis [38, 12] is necessary. Selecting attributes or attribute combinations is as much of an art as a science, and therefore one has to rely on good judgment and on a healthy multidisciplinary debate to progressively reveal missing ones.

## 6.2 Limitations and Future Work

While the advantages of the experimental method are clear, our proposed method does not exempt researchers from exercising attention and good judgment. In particular, while our method greatly reduces unwanted correlations with annotated variables, it does not eliminate them completely, nor does it account for hidden confounders [77], and one will need to keep a sharp eye out for both.

As an example of the first, we found that our method often adds facial hair to male faces when increasing hair length (see Fig. 18). This is likely a reason for why our classifiers did not have higher error rates for males with longer hair (see Fig. 13). As an example of the second, we found that our method tends to synthesize earrings when modifying a dark-skinned face to look female (see Fig. 19). Depending on culture, earrings may or may not be relevant to the definition of gender. If this is an unwanted correlation, one ought to add earrings to the annotation pipeline so that it may be "orthogonalized away" by the synthesis method. Scientists building an industry-grade system for measuring face analysis bias will want to consider including a more exhaustive set of factors. A significant advantage of an approach that is based on synthetic images and human annotation is thus the following: *as soon as one residual correlation is discovered it may be systematically annotated, compensated for in the analysis, and mitigated in the synthesis.*

A number of refinements in face synthesis will make our experimental method more practical and powerful. First, many of the faces we generated contained visible artifacts (see Fig. 8), which we eliminated by human annotation – even subtle artifacts can affect classifier outputs, as revealed by the literature on adversarial examples [75]. Second, we do not yet have tools to estimate the sets of physiognomies and attribute combinations that can and cannot be produced by a given generator. Current GANs are known to have difficulties in generating data outside of their training distributions. Third, we observed a bias of StyleGAN2 towards generating Caucasian faces when sampling from its latent distribution. While our method can compensate for biases through carefully oriented traversals calibrated by human annotations, it would be clearly better to start from unbiased synthesis methods. We are hopeful that these shortcomings will be incrementally resolved by a combination of training sets with increased diversity of attributes like ethnicity, gender, personal style, and age, as well as better models.

Our first-order technique for controlling synthesis can also be improved. A better understanding of the geometry of face space will hopefully yield more accurate global coordinate systems. These, in turn, will help reduce residual biases in synthetic transects, which we currently mitigate by having transects annotated by hand.

Finally, extending our method beyond gender classification to more complex tasks, such as face recognition, is not straightforward in practice and will require further study.

# Acknowledgments

# References

[1] `https://aws.amazon.com/sagemaker/groundtruth/`

[2] Albiero, V., KS, K., Vangara, K., Zhang, K., King, M.C., Bowyer, K.W.: Analysis of gender inequality in face recognition accuracy. In: Proceedings of the IEEE Winter Conference on Applications of Computer Vision Workshops. pp. 81–89 (2020)

[3] Alvi, M., Zisserman, A., Nellåker, C.: Turning a blind eye: Explicit removal of biases and variation from deep neural network embeddings. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 0–0 (2018)

[4] Anderson, D.J., Adolphs, R.: A framework for studying emotions across species. Cell **157**(1), 187–200 (2014)

[5] Angrist, J.D., Imbens, G.W.: Identification and estimation of local average treatment effects. Tech. rep., National Bureau of Economic Research (1995)

[6] Barron, J.L., Fleet, D.J., Beauchemin, S.S.: Performance of optical flow techniques. International journal of computer vision **12**(1), 43–77 (1994)

[7] Bau, D., Zhu, J.Y., Strobelt, H., Zhou, B., Tenenbaum, J.B., Freeman, W.T., Torralba, A.: Gan dissection: Visualizing and understanding generative adversarial networks. In: Proceedings of the International Conference on Learning Representations (ICLR) (2019)

[8] Bertrand, M., Mullainathan, S.: Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination. American economic review **94**(4), 991–1013 (2004)

[9] Bowyer, K., Phillips, P.J.: Empirical evaluation techniques in computer vision. IEEE Computer Society Press (1998)

[10] Brandao, M.: Age and gender bias in pedestrian detection algorithms. arXiv preprint arXiv:1906.10490 (2019)

[11] Buhrmester, M., Kwang, T., Gosling, S.D.: Amazon's mechanical turk: A new source of inexpensive, yet high-quality data? (2016)

[12] Buolamwini, J., Gebru, T.: Gender shades: Intersectional accuracy disparities in commercial gender classification. In: Conference on fairness, accountability and transparency. pp. 77–91 (2018)

[13] Chang, C.H., Creager, E., Goldenberg, A., Duvenaud, D.: Explaining image classifiers by counterfactual generation. In: ICLR (2018)

[14] Corbett-Davies, S., Goel, S.: The measure and mismeasure of fairness: A critical review of fair machine learning. arXiv preprint arXiv:1808.00023 (2018)

[15] Cortes, C., Vapnik, V.: Support-vector networks. Machine learning **20**(3), 273–297 (1995)

[16] Dabkowski, P., Gal, Y.: Real time image saliency for black box classifiers. In: Advances in Neural Information Processing Systems. pp. 6967–6976 (2017)

[17] Darwin, C., Prodger, P.: The expression of the emotions in man and animals. Oxford University Press, USA (1998)

[18] Das, A., Dantcheva, A., Bremond, F.: Mitigating bias in gender, age and ethnicity classification: a multi-task convolution neural network approach. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 0–0 (2018)

[19] Denton, E., Hutchinson, B., Mitchell, M., Gebru, T.: Detecting bias with generative counterfactual face attribute augmentation. arXiv preprint arXiv:1906.06439 (2019)

[20] Drozdowski, P., Rathgeb, C., Dantcheva, A., Damer, N., Busch, C.: Demographic bias in biometrics: A survey on an emerging challenge. IEEE Transactions on Technology and Society (2020)

[21] Egan, S.K., Perry, D.G.: Gender identity: a multidimensional analysis with implications for psychosocial adjustment. Developmental psychology **37**(4), 451 (2001)

[22] Fei-Fei, L., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In: 2004 conference on computer vision and pattern recognition workshop. pp. 178–178. IEEE (2004)

[23] Fitzpatrick, T.B.: The validity and practicality of sun-reactive skin types i through vi. Archives of dermatology **124**(6), 869–871 (1988)

[24] Fong, R.C., Vedaldi, A.: Interpretable explanations of black boxes by meaningful perturbation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3429–3437 (2017)

[25] Gelman, A., Hill, J.: Data analysis using regression and multilevel/hierarchical models. Cambridge university press (2006)

[26] Goyal, Y., Wu, Z., Ernst, J., Batra, D., Parikh, D., Lee, S.: Counterfactual visual explanations. In: Proceedings of the 36th International Conference on Machine Learning. vol. 97, pp. 2376–2384. PMLR (2019)

[27] Grother, P., Ngan, M., Hanaoka, K.: Ongoing face recognition vendor test (frvt) part 1: Verification. National Institute of Standards and Technology, Tech. Rep (2018)

[28] Grother, P.J., Ngan, M.L., Hanaoka, K.K.: Ongoing face recognition vendor test (frvt) part 2: identification. Tech. rep. (2018)

[29] Hanaoka, P.G.N.K.: Face recognition vendor test (frvt)part 3: Demographic effects. IR 8280, NIST, https://doi.org/10.6028/NIST.IR.8280 (2019)

[30] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

[31] Hébert-Johnson, U., Kim, M.P., Reingold, O., Rothblum, G.N.: Calibration for the (computationally-identifiable) masses. arXiv preprint arXiv:1711.08513 (2017)

[32] Heckman, J.J., Vytlacil, E.J.: Instrumental variables, selection models, and tight bounds on the average treatment effect. In: Econometric Evaluation of Labour Market Policies, pp. 1–15. Springer (2001)

[33] Hendricks, L.A., Burns, K., Saenko, K., Darrell, T., Rohrbach, A.: Women also snowboard: Overcoming bias in captioning models. In: European Conference on Computer Vision. pp. 793–811. Springer (2018)

[34] Hoerl, A.E., Kennard, R.W.: Ridge regression: Biased estimation for nonorthogonal problems. Technometrics **12**(1), 55–67 (1970)

[35] Kärkkäinen, K., Joo, J.: Fairface: Face attribute dataset for balanced race, gender, and age. arXiv preprint arXiv:1908.04913 (2019)

[36] Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4401–4410 (2019)

[37] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. arXiv preprint arXiv:1912.04958 (2019)

[38] Kearns, M., Neel, S., Roth, A., Wu, Z.S.: Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. arXiv preprint arXiv:1711.05144 (2017)

[39] Kearns, M., Neel, S., Roth, A., Wu, Z.S.: An empirical study of rich subgroup fairness for machine learning. In: Proceedings of the Conference on Fairness, Accountability, and Transparency. pp. 100–109 (2019)

[40] Kearns, M., Roth, A.: The Ethical Algorithm: The Science of Socially Aware Algorithm Design. Oxford University Press (2019)

[41] Khosla, A., Zhou, T., Malisiewicz, T., Efros, A.A., Torralba, A.: Undoing the damage of dataset bias. In: European Conference on Computer Vision. pp. 158–171. Springer (2012)

[42] Kim, B., M., W., Gilmer, J., C., C., J., W., , Viegas, F., Sayres, R.: Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV) . ICML (2018)

[43] Klare, B.F., Burge, M.J., Klontz, J.C., Bruegge, R.W.V., Jain, A.K.: Face recognition performance: Role of demographic information. IEEE Transactions on Information Forensics and Security **7**(6), 1789–1801 (2012)

[44] Kleinberg, J., Ludwig, J., Mullainathany, S., Sunstein, C.R.: Discrimination in the age of algorithms. Published by Oxford University Press on behalf of The John M. Olin Center for Law, Economics and Business at Harvard Law School (2019), https://academic.oup.com/jla/article-abstract/doi/10.1093/jla/laz001/5476086

[45] Kortylewski, A., Egger, B., Schneider, A., Gerig, T., Morel-Forster, A., Vetter, T.: Empirically analyzing the effect of dataset biases on deep face recognition systems. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. pp. 2093–2102 (2018)

[46] 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: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. pp. 0–0 (2019)

[47] Krishnapriya, K.S., Vangara, K., King, M., Albiero, V., Bowyer, K.: Characterizing the variability in face recognition accuracy relative to race. ArXiv 1904.07325 (4 2019)

[48] Li, Y., Vasconcelos, N.: Repair: Removing representation bias by dataset resampling. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9572–9581 (2019)

[49] Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: Proceedings of International Conference on Computer Vision (ICCV) (2015)

[50] Lohr, S.: Facial recognition is accurate, if you're a white guy. New York Times (February 9 2018), https://nyti.ms/2BNurVq

[51] Lu, B., Chen, J.C., Castillo, C.D., Chellappa, R.: An experimental evaluation of covariates effects on unconstrained face verification. IEEE Transactions on Biometrics, Behavior, and Identity Science **1**(1), 42–55 (2019)

[52] Mendez, M.A.: Presumptions of discriminatory motive in title vii disparate treatment cases. Stanford Law Review pp. 1129–1162 (1980)

[53] Merkatz, R.B., Temple, R., Sobel, S., Feiden, K., Kessler, D.A., on Women in Clinical Trials, W.G.: Women in clinical trials of new drugs–a change in food and drug administration policy. New England Journal of Medicine **329**(4), 292–296 (1993)

[54] Merler, M., Ratha, N., Feris, R.S., Smith, J.R.: Diversity in faces. arXiv preprint arXiv:1901.10436 (2019)

[55] Mullainathan, S.: Biased algorithms are easier to fix than biased people. New York Times (6 Dec 2019)

[56] Muthukumar, V., Pedapati, T., Ratha, N., Sattigeri, P., Wu, C.W., Kingsbury, B., Kumar, A., Thomas, S., Mojsilovic, A., Varshney, K.R.: Understanding unequal gender classification accuracy from face images. arXiv preprint arXiv:1812.00099 (2018)

[57] Oreopoulos, P.: Estimating average and local average treatment effects of education when compulsory schooling laws really matter. American Economic Review **96**(1), 152–175 (2006)

[58] Pearl, J.: Causality. Cambridge university press (2009)

[59] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12**, 2825–2830 (2011)

[60] Phillips, P.J., Grother, P., Micheals, R., Blackburn, D.M., Tabassi, E., Bone, M.: Face recognition vendor test 2002. In: 2003 IEEE International SOI Conference. Proceedings (Cat. No. 03CH37443). p. 44. IEEE (2003)

[61] Phillips, P.J., Wechsler, H., Huang, J., Rauss, P.J.: The feret database and evaluation procedure for face-recognition algorithms. Image and vision computing **16**(5), 295–306 (1998)

[62] Phillips, P.J., Yates, A.N., Hu, Y., Hahn, C.A., Noyes, E., Jackson, K., Cavazos, J.G., Jeckeln, G., Ranjan, R., Sankaranarayanan, S., et al.: Face recognition accuracy of forensic examiners, superrecognizers, and face recognition algorithms. Proceedings of the National Academy of Sciences **115**(24), 6171–6176 (2018)

[63] Pocock, S.J., Assmann, S.E., Enos, L.E., Kasten, L.E.: Subgroup analysis, covariate adjustment and baseline comparisons in clinical trial reporting: current practiceand problems. Statistics in medicine **21**(19), 2917–2930 (2002)

[64] Ponce, J., Berg, T.L., Everingham, M., Forsyth, D.A., Hebert, M., Lazebnik, S., Marszalek, M., Schmid, C., Russell, B.C., Torralba, A., et al.: Dataset issues in object recognition. In: Toward category-level object recognition, pp. 29–48. Springer (2006)

[65] Robinson, L.D., Jewell, N.P.: Some surprising results about covariate adjustment in logistic regression models. International Statistical Review/Revue Internationale de Statistique pp. 227–240 (1991)

[66] Rubin, D.B.: Matched sampling for causal effects. Cambridge University Press (2006)

[67] Rutherglen, G.: Disparate impact under title vii: an objective theory of discrimination. Virginia Law Review pp. 1297–1345 (1987)

[68] Ryu, H.J., Adam, H., Mitchell, M.: Inclusivefacenet: Improving face attribute detection with race and gender diversity. arXiv preprint arXiv:1712.00193 (2017)

[69] Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 618–626 (2017)

[70] Shen, Y., Gu, J., Tang, X., Zhou, B.: Interpreting the latent space of gans for semantic face editing. arXiv preprint arXiv:1907.10786 (2019)

[71] Simon, V.: Wanted: women in clinical trials (2005)

[72] Simonyan, K., Vedaldi, A., Zisserman, A.: Deep inside convolutional networks: Visualising image classification models and saliency maps (2013)

[73] Singla, S., Pollack, B., Chen, J., Batmanghelich, K.: Explanation by progressive exaggeration. arXiv preprint arXiv:1911.00483 (2019)

[74] Sundararajan, M., Taly, A., Yan, Q.: Axiomatic attribution for deep networks. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 3319–3328. JMLR. org (2017)

[75] Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)

[76] Torralba, A., Efros, A.A., et al.: Unbiased look at dataset bias. In: CVPR. vol. 1, p. 7 (2011)

[77] VanderWeele, T.J., Shpitser, I.: On the definition of a confounder. Annals of statistics **41**(1), 196 (2013)

[78] Willan, A.R., Briggs, A.H., Hoch, J.S.: Regression methods for covariate adjustment and subgroup analysis for non-censored cost-effectiveness data. Health economics **13**(5), 461–475 (2004)

[79] Wilson, E.B.: Probable inference, the law of succession, and statistical inference. Journal of the American Statistical Association **22**(158), 209–212 (1927)

[80] Xiao, T., Hong, J., Ma, J.: Elegant: Exchanging latent encodings with gan for transferring multiple face attributes. In: Proceedings of the European conference on computer vision (ECCV). pp. 168–184 (2018)

[81] Youla, D.: Mathematical theory of image restoration by the method of convex projections. Image recovery: theory and application pp. 29–77 (1987)

[82] Zhou, B., Bau, D., Oliva, A., Torralba, A.: Interpreting deep visual representations via network dissection. IEEE transactions on pattern analysis and machine intelligence (2018)