

Responsible AI: Measuring accuracy and bias in Computer Vision Systems

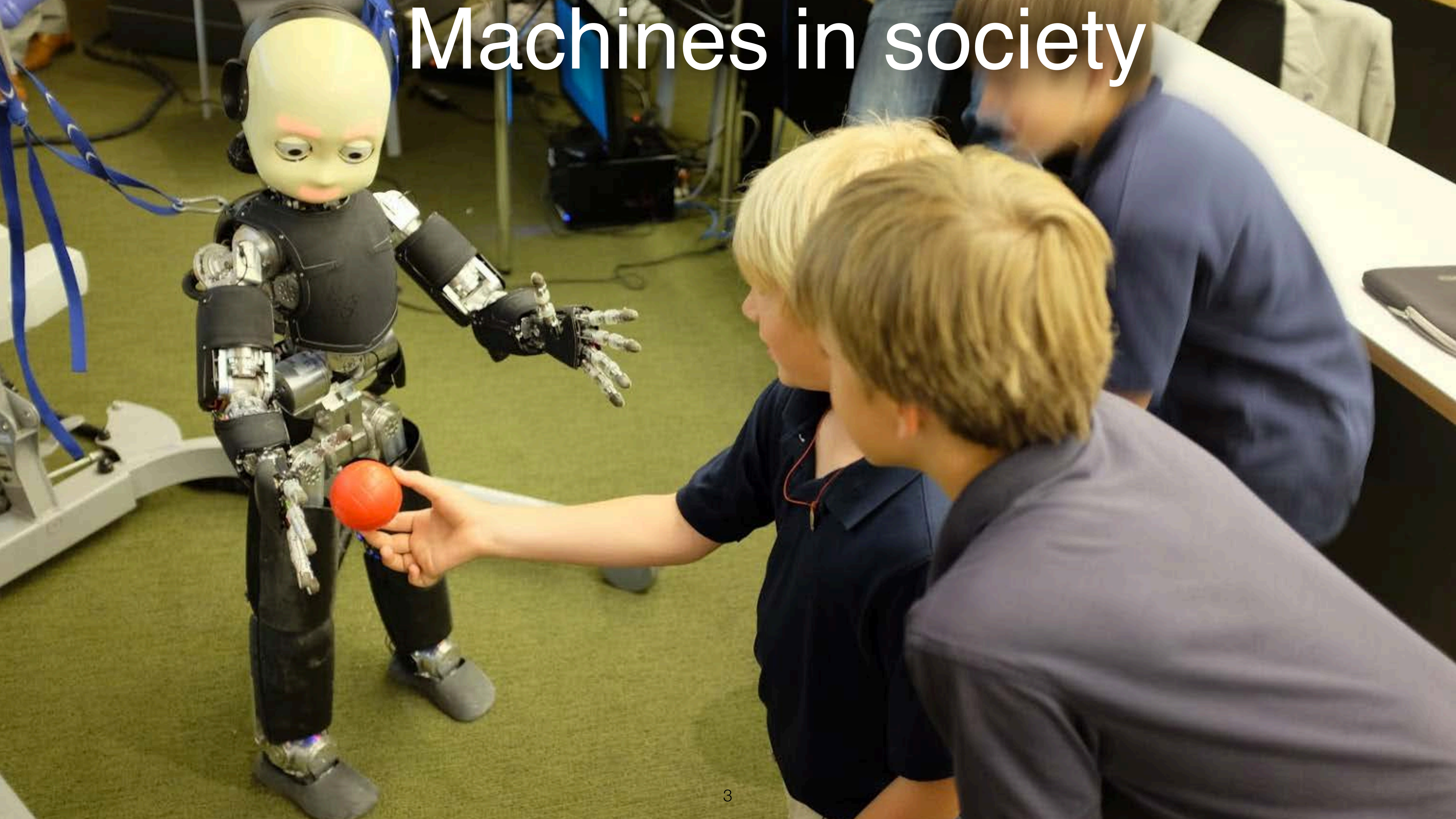
Pietro Perona
Caltech and AWS

Ellis Summer School on Large Scale AI
Modena - 18-22 September 2023

Plan

- Responsible AI
- Bias and fairness
- Correlation vs causation
- Bias in face analysis and recognition
- Bias in generative models
- Fixing bias
- Conclusions

Machines in society



Responsible AI

- Fairness, bias
- Robustness
- Explainability
- Safety
- Privacy
- IP



WORLD & NATION

Identical twins. Identical asylum claims. Very different luck at the border

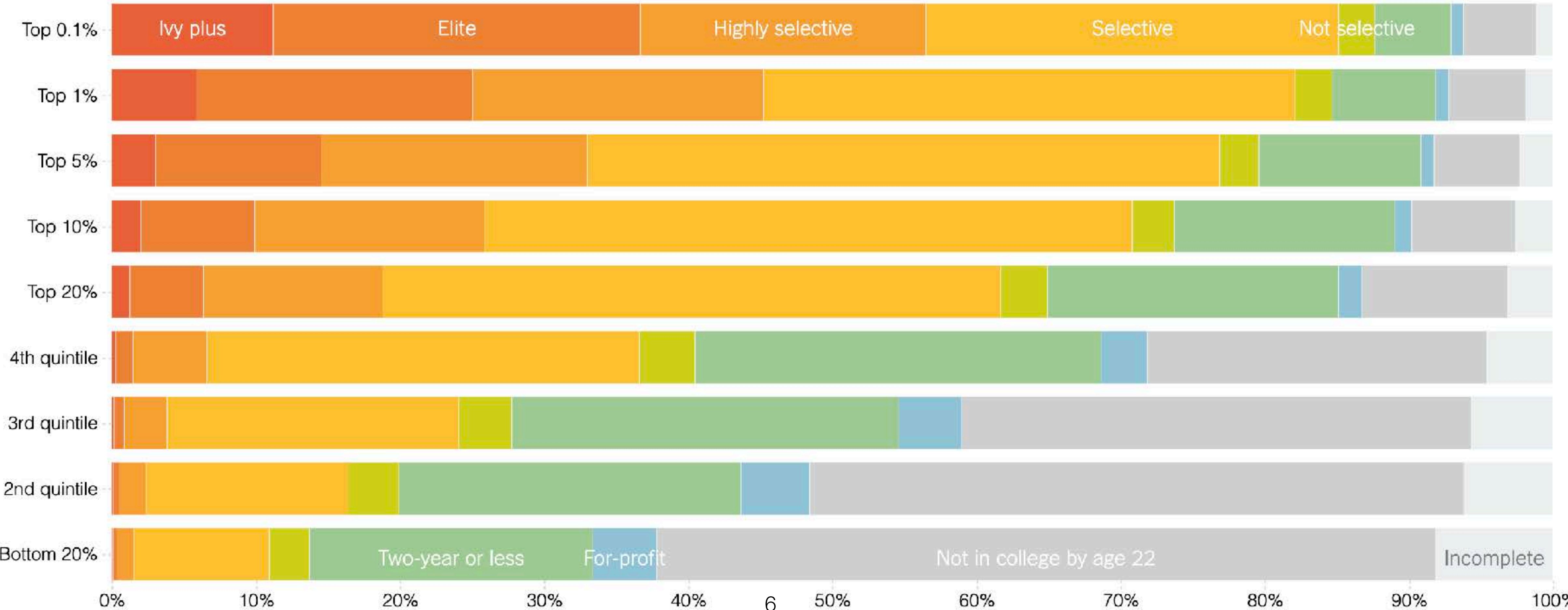


Save

Q

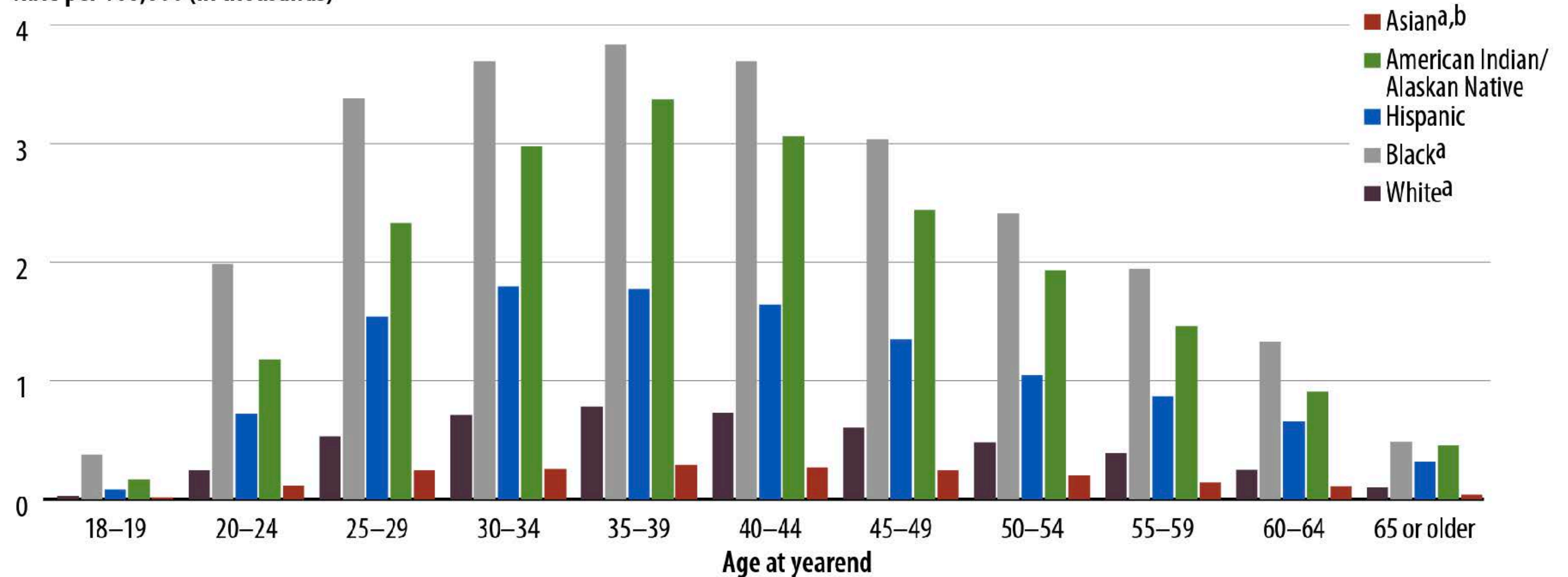
Where today's 25-year-olds went to college, grouped by their parents' income

About four in 10 students from the top 0.1 percent attend an Ivy League or elite university, roughly equivalent to the share of students from poor families who attend any two- or four-year college.



Imprisonment rates of male U.S. residents, based on sentenced prisoners under the jurisdiction of state or federal correctional authorities, by demographic characteristics, December 31, 2021

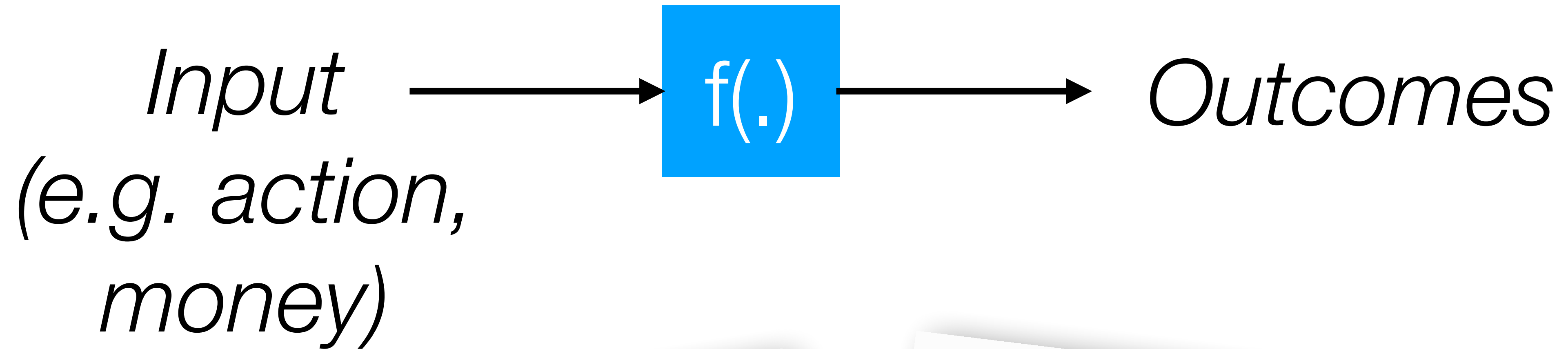
Rate per 100,000 (in thousands)



Males: 1,084,717

Females: 78,948

Fairness



AEA Papers and Proceedings 2018, 108: 22–27
<https://doi.org/10.1257/pandp.20181018>

ADVANCES IN BIG DATA RESEARCH IN ECONOMICS

Algorithmic Fairness[†]

By JON KLEINBERG, JENS LUDWIG, SENDHIL MULLAINATHAN, AND ASHESH RAMBACHAN*

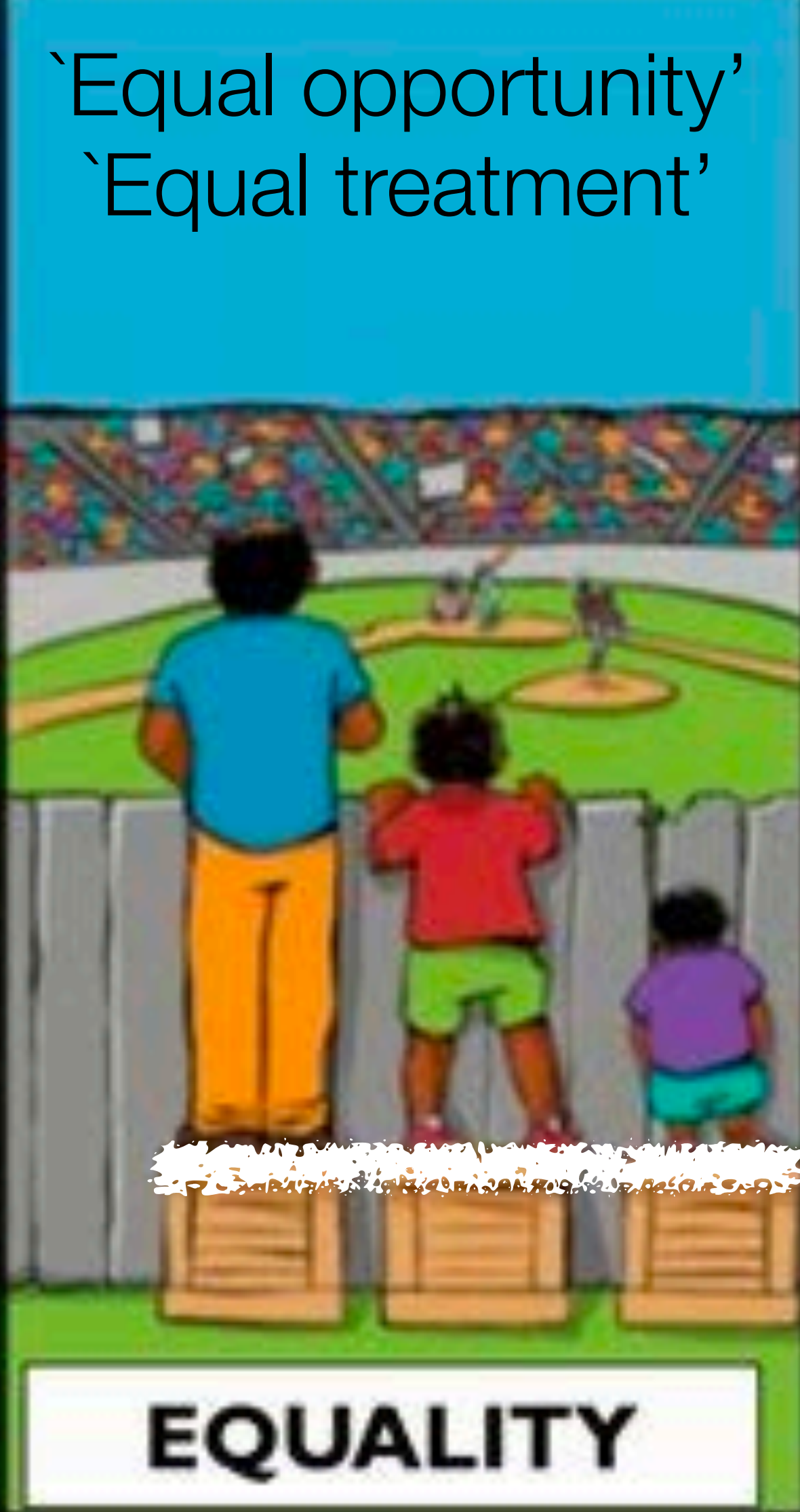
Annual Review of Statistics and Its Application
Algorithmic Fairness: Choices,
Assumptions, and Definitions

Shira Mitchell,¹ Eric Potash,² Solon Barocas,^{3,4}
Alexander D'Amour,⁵ and Kristian Lum⁶



REALITY

One gets **more than** is needed, while the other gets **less than** is needed. Thus, a huge disparity is created.



EQUALITY

The assumption is that **everyone benefits from the same supports**. This is considered to be equal treatment.



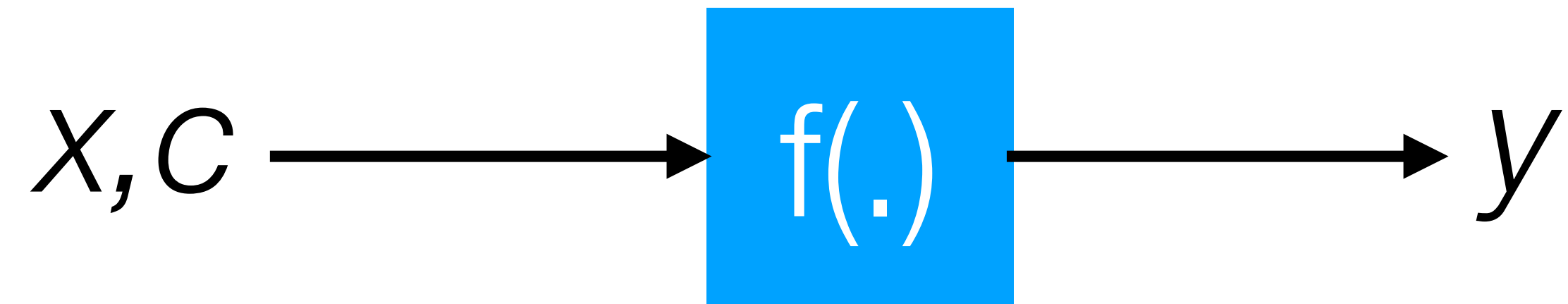
EQUITY

Everyone gets the support they need, which produces equity.

‘Equal opportunity’
‘Equal treatment’

‘Equal outcomes’

In equations



$x = \text{action, state}$

$c = \text{protected attribute}$

$y = \text{reward, outcome}$

$$p(y = 1|x, c) = p(y = 1|x)$$

$$p(y = 1|x, c = 0) = p(y = 1|x, c = 1)$$

Examples

New York City marathon

Men

2010	Gebregziabher Gebremariam	 Ethiopia	2:08:14
2011	Geoffrey Mutai	 Kenya	2:05:05
2012	Cancelled due to Hurricane Sandy		
2013	Geoffrey Mutai (2)	 Kenya	2:08:24
2014	Wilson Kipsang	 Kenya	2:10:59
2015	Stanley Biwott	 Kenya	2:10:34
2016	Ghirmay Ghebreslassie	 Eritrea	2:07:51
2017	Geoffrey Kamworor	 Kenya	2:10:53
2018	Lelisa Desisa	 Ethiopia	2:05:59
2019	Geoffrey Kamworor (2)	 Kenya	2:08:13

Women

2010	Edna Kiplagat	 Kenya	2:28:20
2011	Firehiwot Dado	 Ethiopia	2:23:15
2012	Cancelled due to Hurricane Sandy		
2013	Priscah Jeptoo	 Kenya	2:25:07
2014	Mary Keitany	 Kenya	2:25:07
2015	Mary Keitany	 Kenya	2:24:25
2016	Mary Keitany	 Kenya	2:24:26
2017	Shalane Flanagan	 United States	2:26:53
2018	Mary Keitany	 Kenya	2:22:48
2019	Joyciline Jepkosgei	 Kenya	2:22:38

Jocelyn Jepkosgei and
Geoffrey Kamworor,
winners of the 2019
NYC marathon



Dressage

Games	Gold	Silver	Bronze
1980 Moscow details	 Elisabeth Theurer on <i>Mane</i> (AUT)	 Yuri Kovshov on <i>Igrok</i> (URS)	 Viktor Ugryumov on <i>St...</i> (URS)
1984 Los Angeles details	 Reiner Klimke on <i>Ahlerich</i> (FRG)	 Anne Grethe Jensen on <i>Marzog</i> (DEN)	 Otto Hofer on <i>Limandus</i> (SUI)
1988 Seoul details	 Nicole Uphoff on <i>Rembrandt</i> (FRG)	 Margit Otto-Crépin on <i>Corlandus</i> (FRA)	 Christine Stuckelberger on <i>...</i> (SUI)
1992 Barcelona details	 Nicole Uphoff on <i>Rembrandt</i> (GER)	 Isabell Werth on <i>Gigolo</i> (GER)	 Klaus Balkenhol on <i>Gold...</i> (GER)
1996 Atlanta details	 Isabell Werth on <i>Gigolo</i> (GER)	 Anky van Grunsven on <i>Bonfire</i> (NED)	 Sven Rothenberger on <i>Weyden</i> (NED)
2000 Sydney details	 Anky van Grunsven on <i>Bonfire</i> (NED)	 Isabell Werth on <i>Gigolo</i> (GER)	 Ulla Salzgeber on <i>Rusty</i> (GER)
2004 Athens details	 Anky van Grunsven on <i>Salinero</i> (NED)	 Ulla Salzgeber on <i>Rusty</i> (GER)	 Beatriz Ferrer-Salat on <i>Beauvalais</i> (ESP)
2008 Beijing details	 Anky van Grunsven on <i>Salinero</i> (NED)	 Isabell Werth on <i>Satchmo</i> (GER)	 Heike Kemmer on <i>Bonaparte</i> (GER)
2012 London details	 Charlotte Dujardin on <i>Valegro</i> (GBR)	 Adelinde Cornelissen on <i>Parzival</i> (NED)	 Laura Bechtolsheimer on <i>Mistral Højris</i> (GBR)
2016 Rio details	 Charlotte Dujardin on <i>Valegro</i> (GBR)	 Isabell Werth on <i>Weihegold</i> (GER)	 Kristina Bröring-Sprehe on <i>Desperados</i> (GER)

Jumping

Games	Gold	Silver	Bronze
1980 Moscow details	 Jan Kowalczyk on <i>Artemor</i> (POL)	 Nikolai Korolkov on <i>Espadron</i> (URS)	 Joaquín Perez Heras on <i>Alymony</i> (MEX)
1984 Los Angeles details	 Joseph Fargis on <i>Touch of Class</i> (USA)	 Conrad Homfeld on <i>Abdullah</i> (USA)	 Heidi Robbiani on <i>Jessica V</i> (SUI)
1988 Seoul details	 Pierre Durand Jr. on <i>Jappeloup</i> (FRA)	 Greg Best on <i>Gem Twist</i> (USA)	 Karsten Huk on <i>Nepomuk</i> (FRG)
1992 Barcelona details	 Ludger Beerbaum on <i>Classic Touch</i> (GER)	 Piet Raymakers on <i>Ratina Z</i> (NED)	 Norman Dello Joio on <i>Iris</i> (USA)
1996 Atlanta details	 Ulrich Kirchhoff on <i>Jus de Pommes</i> (GER)	 Willi Melliger on <i>Calvaro V</i> (SUI)	 Alexandra Ledermann on <i>Rochet M</i> (FRA)
2000 Sydney details	 Jeroen Dubbeldam on <i>De Sjiem</i> (NED)	 Albert Voorn on <i>Lando</i> (NED)	 Khaled Al-Eid on <i>Khashm Al Aan</i> (KSA)
2004 Athens details	 Rodrigo Pessoa on <i>Baloubet du Rouet</i> (BRA)	 Chris Kappler on <i>Royal Kaliber</i> (USA)	 Marco Kutscher on <i>Montender</i> (GER)
2008 Beijing details	 Eric Lamaze on <i>Hickstead</i> (CAN)	 Rolf-Göran Bengtsson on <i>Ninja</i> (SWE)	 Beezie Madden on <i>Authentic</i> (USA)
2012 London details	 Steve Guerdat on <i>Nino Des Buissonets</i> (SUI)	 Gerco Schroder on <i>London</i> (NED)	 Cian O'Connor on <i>Blue Loyd</i> (IRL)
2016 Rio details	 Nick Skelton on <i>Big Star</i> (GBR)	 Peder Fredricson on <i>All In</i> (SWE)	 Eric Lamaze on <i>Fine Lady</i> (CAN)



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

Machine Bias


There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Matu and Lauren Kirchner, ProPublica

May 23, 2016

Parole

Department of Justice
UNITED STATES PENITENTIARY
Atlanta, Georgia



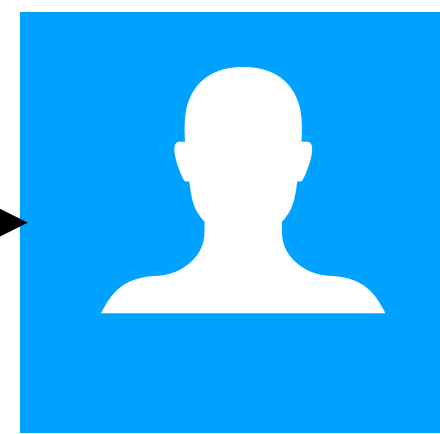
Received
From *Ill. Chicago*
Crime *Via Income Tax Law*
Sentence: *10* yrs. — mos. — days
Date of sentence *Oct 24 - 1931*
Sentence begins *May 4 - 1932*
Sentence expires *May 3 - 1942*
Good time sentence expires *Jan 19 - 1939*
Date of birth *1-17-99* Occupation *Gambler*
Birthplace *N.Y.* Nationality
Age *33* Comp *Med fr*
Height *5'-10 1/2"* Eyes *Grey*
Weight *255* Hair *dk. brn*
Build *stout*
Residence *Chicago, Ill*

Scars and marks: *b. sc. 4" across cheek 2" on L ear - Vert sc 2 1/2" on L jaw - b. sc 2 1/2" - 2" on L ear on neck*

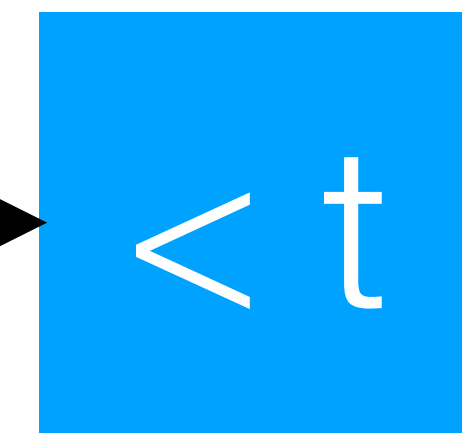
CRIMINAL HISTORY

NAME	NUMBER	CITY	DATE	CHARGE	DISPOSITION
<i>My City</i>			<i>1919</i>	<i>Dis Cond</i>	<i>Discharged</i>
<i>Chicago Ill</i>			<i>1923</i>	<i>Traffic Vio</i>	<i>Dismissed</i>
<i>Do</i>			<i>5-8-24</i>	<i>Murder Witness</i>	<i>Released</i>
<i>Do</i>			<i>6-7-26</i>	<i>Vio U.P.A.</i>	<i>Dismissed</i>
<i>Do</i>			<i>7-28-26</i>	<i>Murder</i>	<i>Charge Withdrawn</i>
<i>Do</i>			<i>10-1-26</i>	<i>Vio U.P.A.</i>	<i>Dismissed</i>
<i>Do</i>			<i>11-12-27</i>	<i>Refused to test</i>	<i>Do</i>
<i>Joliet Ill</i>			<i>12-22-27</i>	<i>Con. Work</i>	<i>Fined \$26.00.00</i>
<i>Phila. Pa</i>			<i>5-11-29</i>	<i>Con. Work</i>	<i>Served 12 mos</i>
<i>Miami Fla</i>			<i>1925</i>	<i>Susp.</i>	<i>Released</i>
<i>Do</i>			<i>5-8-30</i>	<i>Do</i>	<i>Held nights release</i>

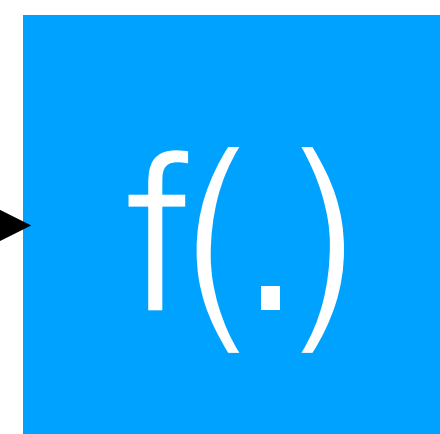
form j n g 55-12-9-31 5-M *for other arrests see Declaration*



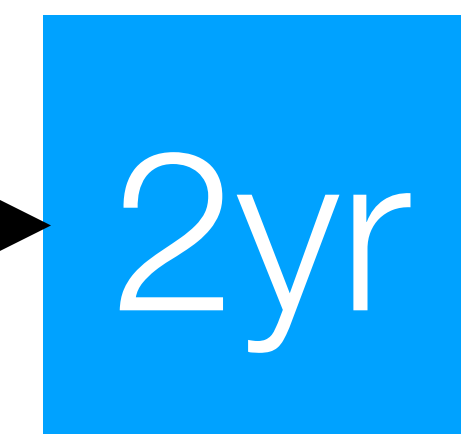
Probability of reoffending



Free



*Risk
1-10*



*Reoffend?
Y/N*



Machine Bias

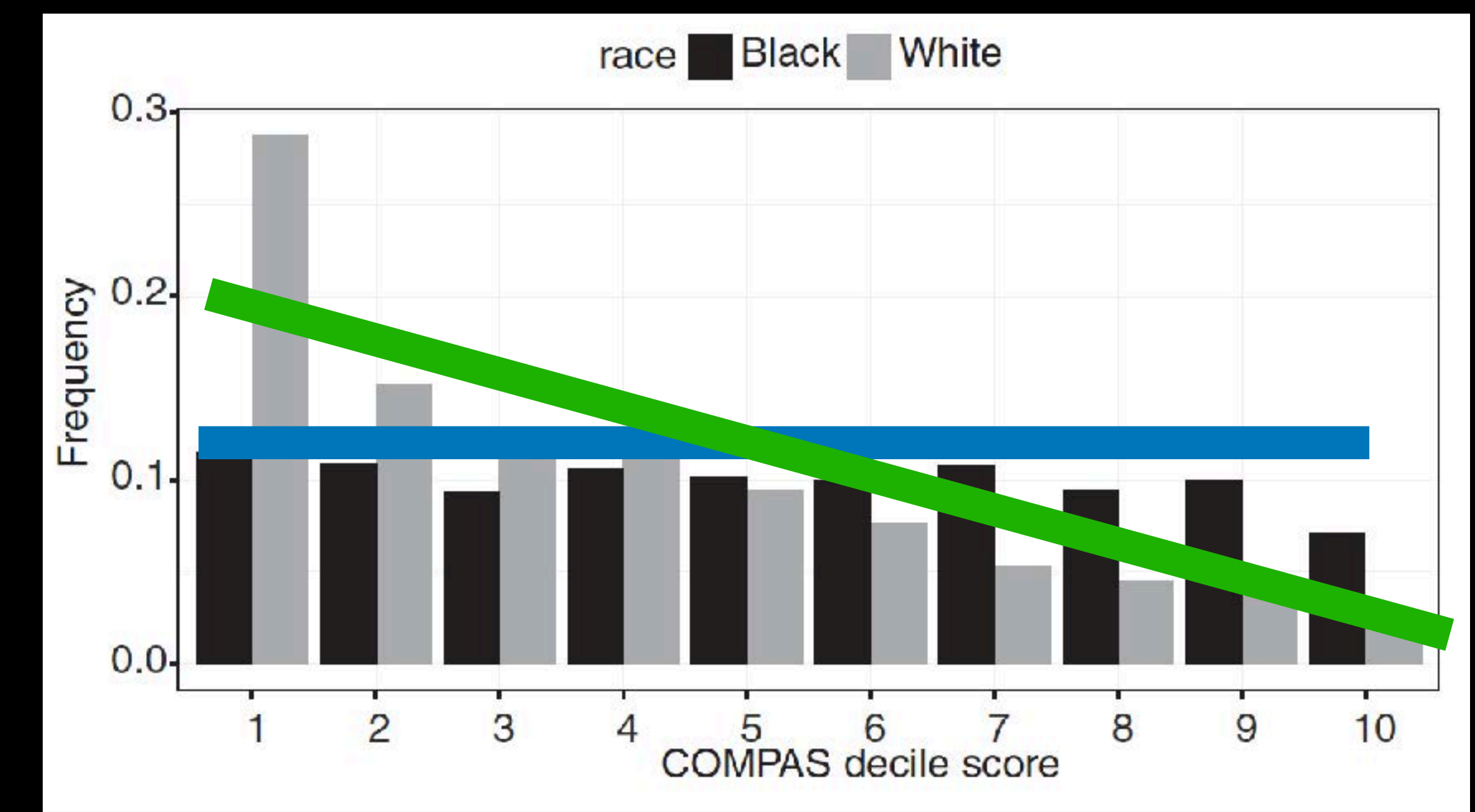
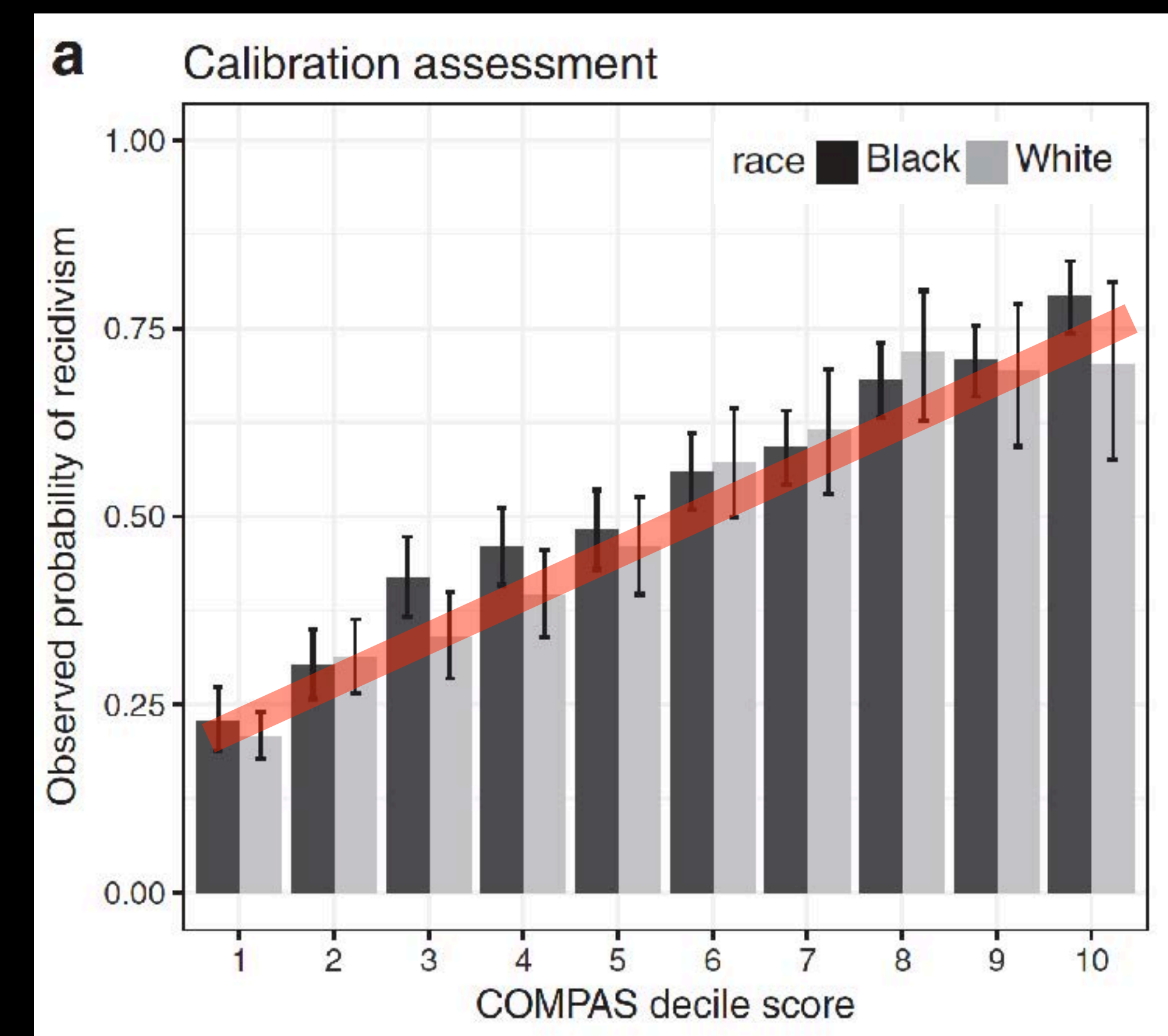
There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

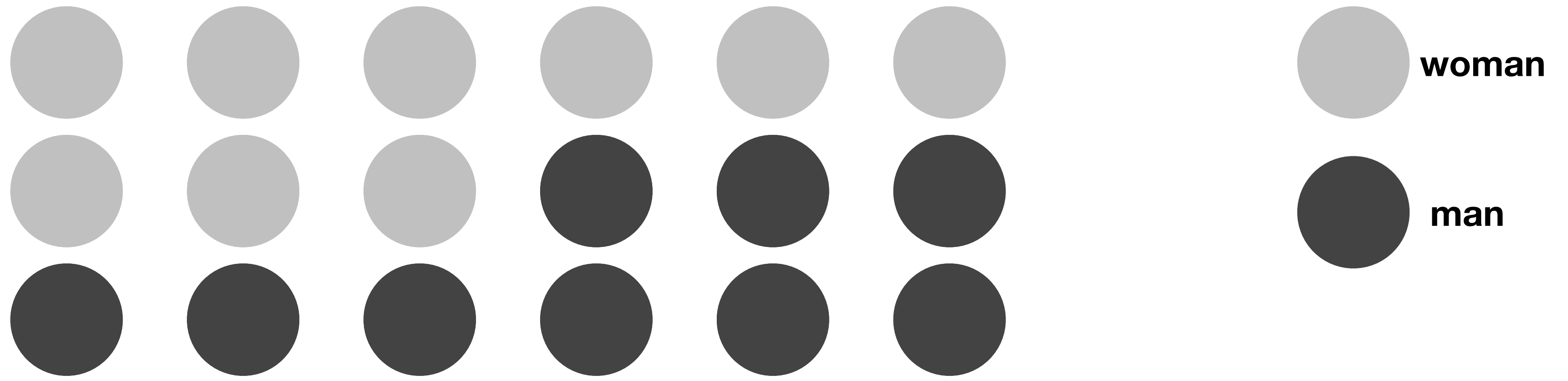
Prediction Fails Differently for Black Defendants

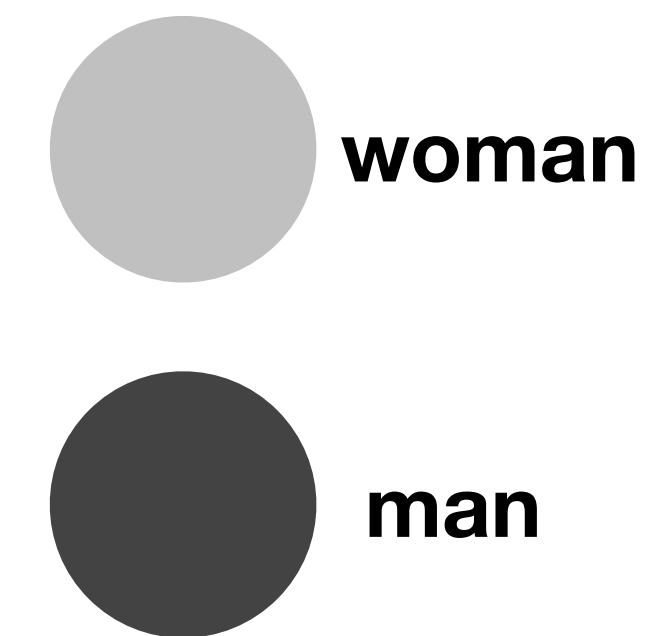
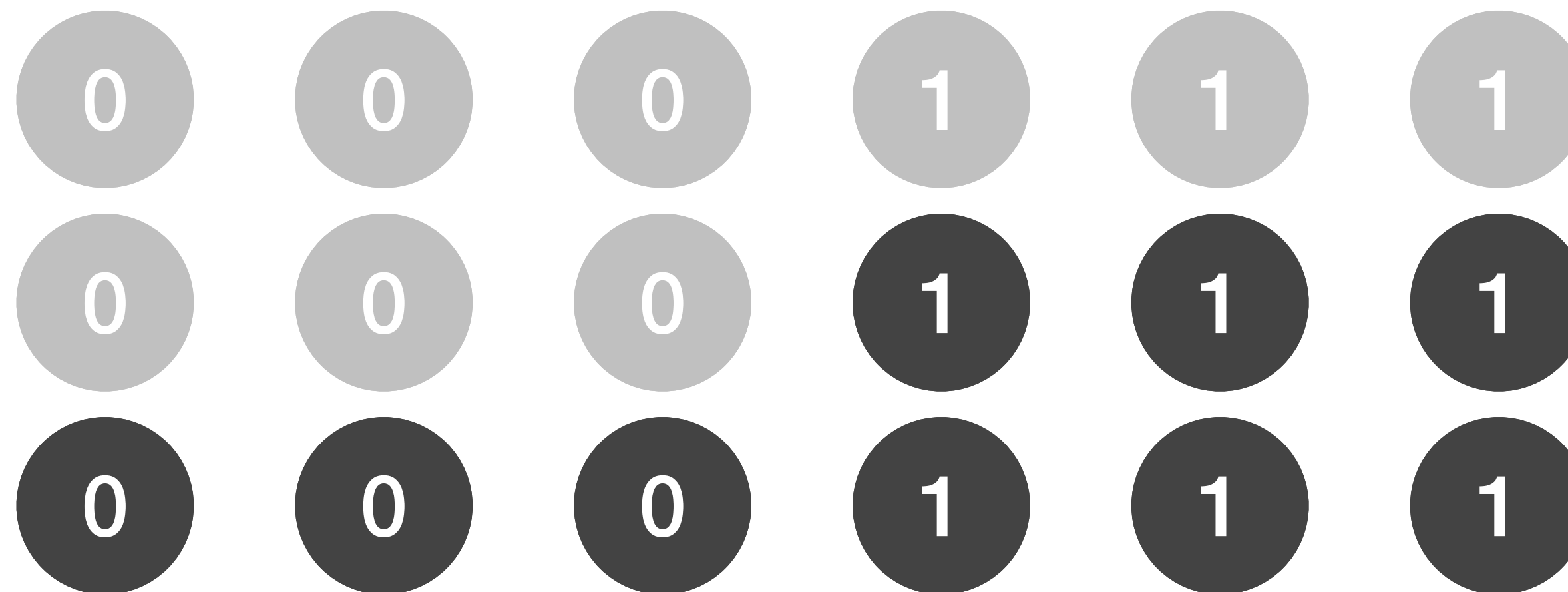
[Pro-Publica 2015]

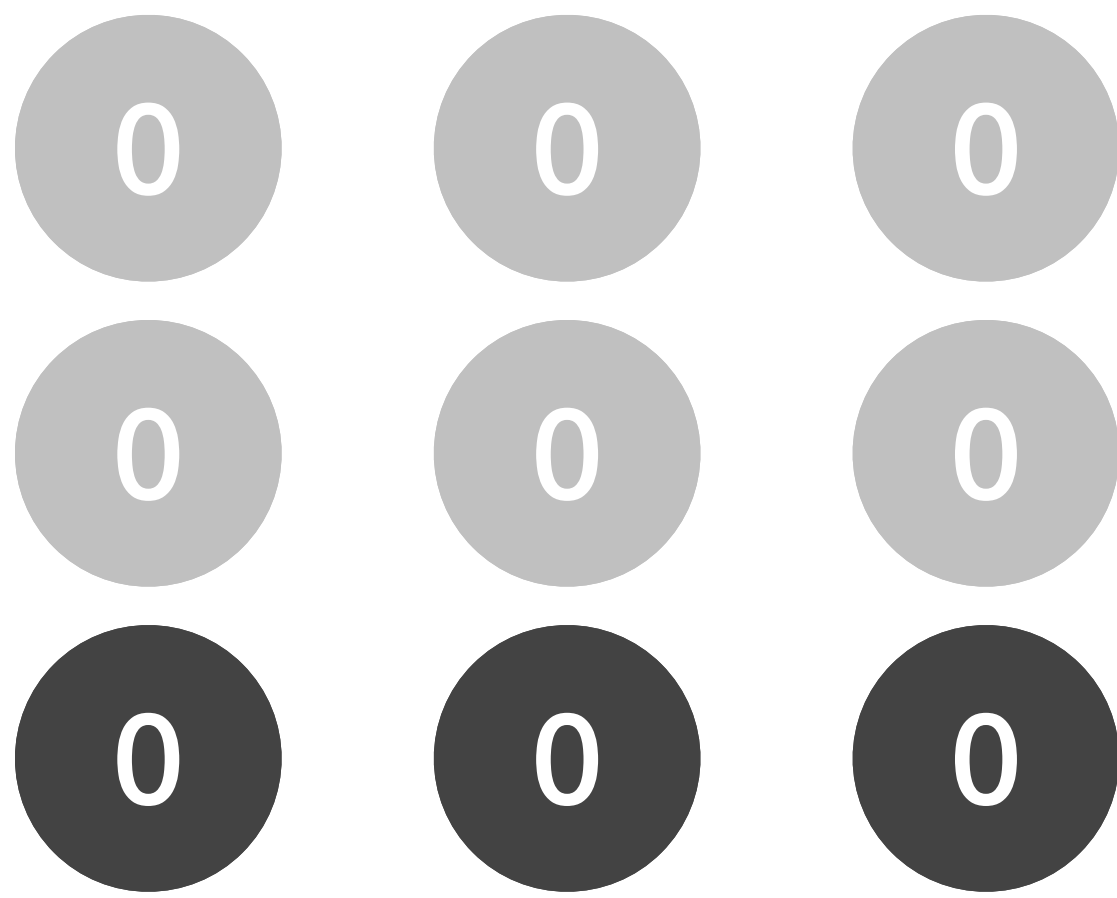
	WHITE	AFRICAN AMERICAN
Non-reoffenders - % labeled hi risk	23.5%	44.9%
Reoffenders - % labeled lower risk	47.7%	28.0%



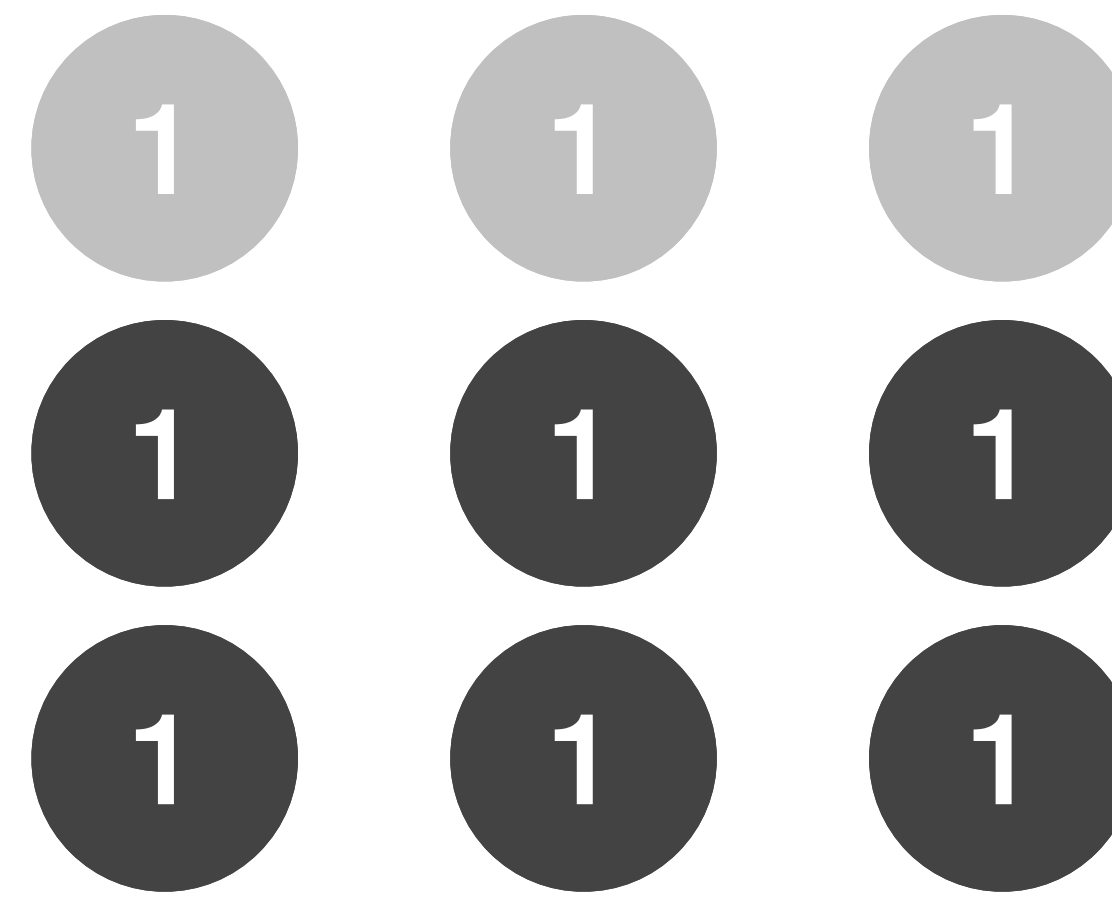
[Chouldechova 2017]



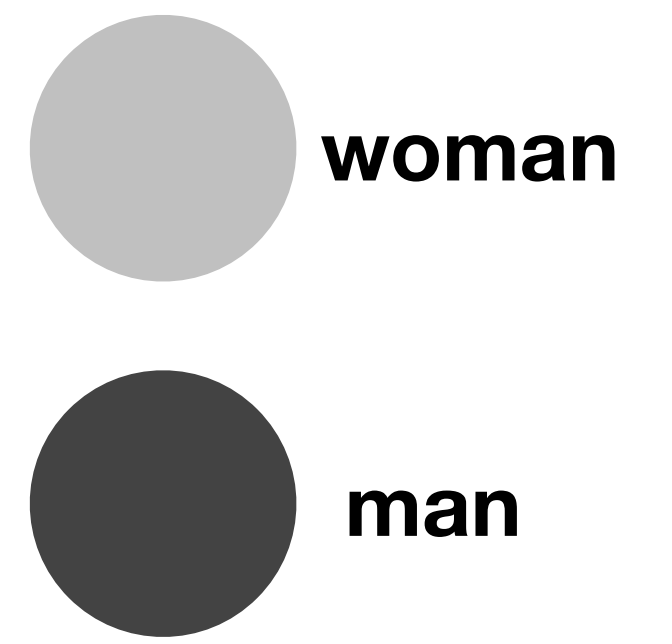


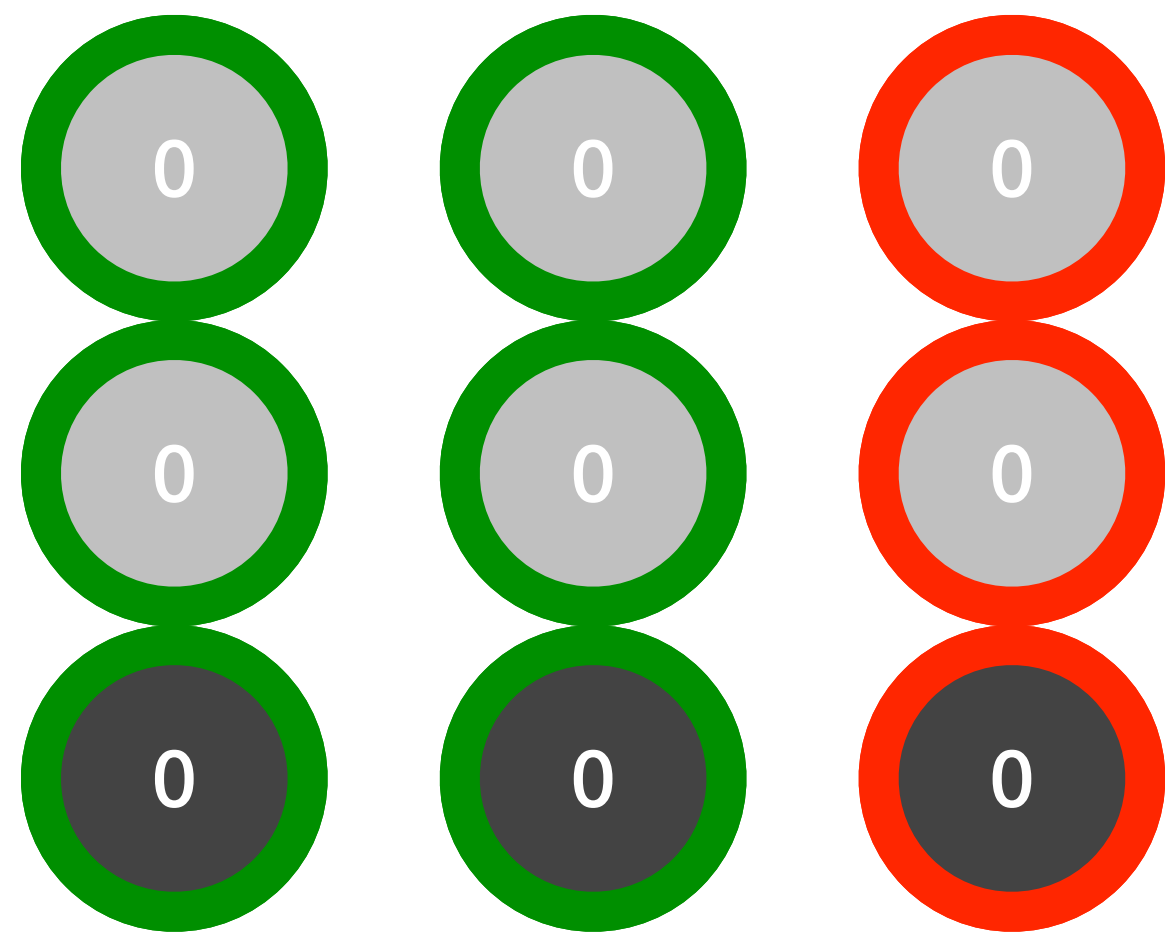


LOW SCORE

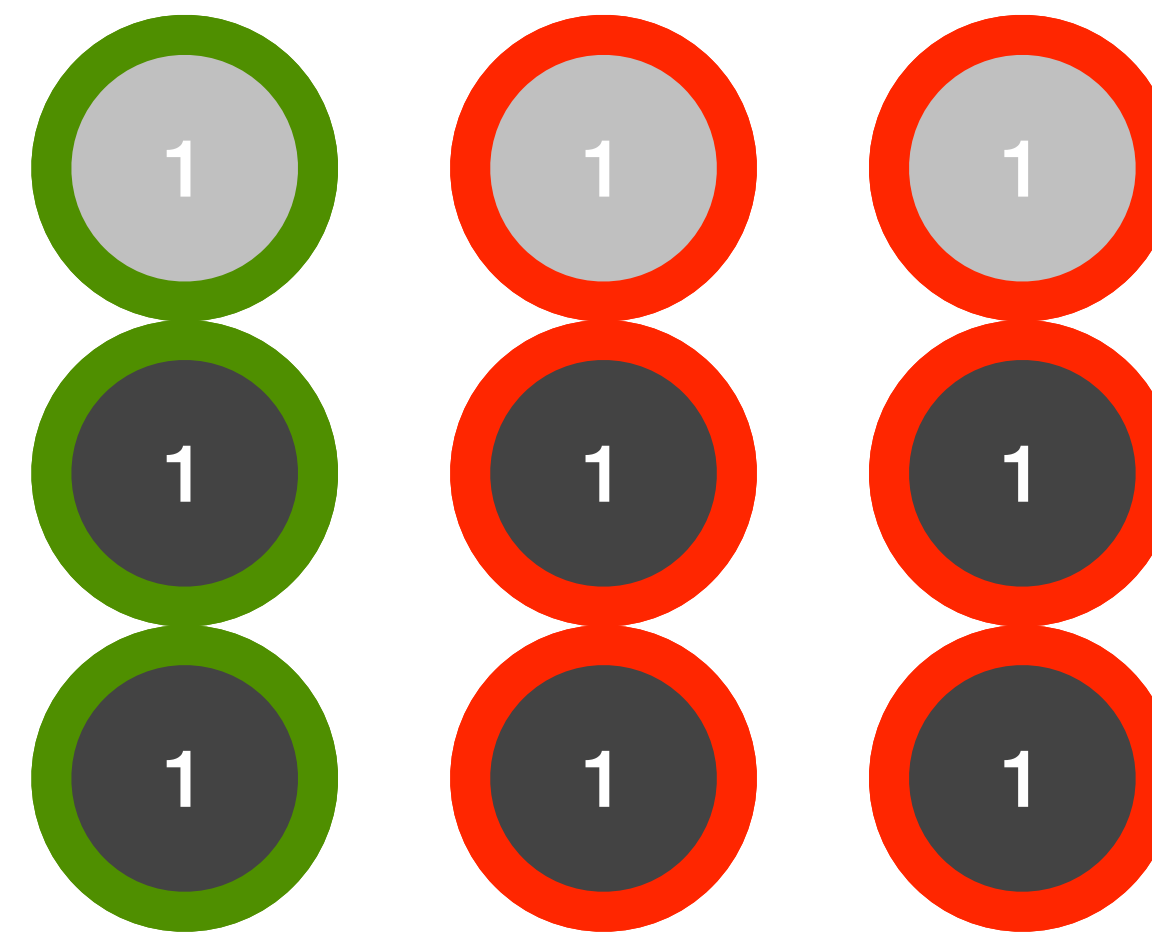


HIGH SCORE

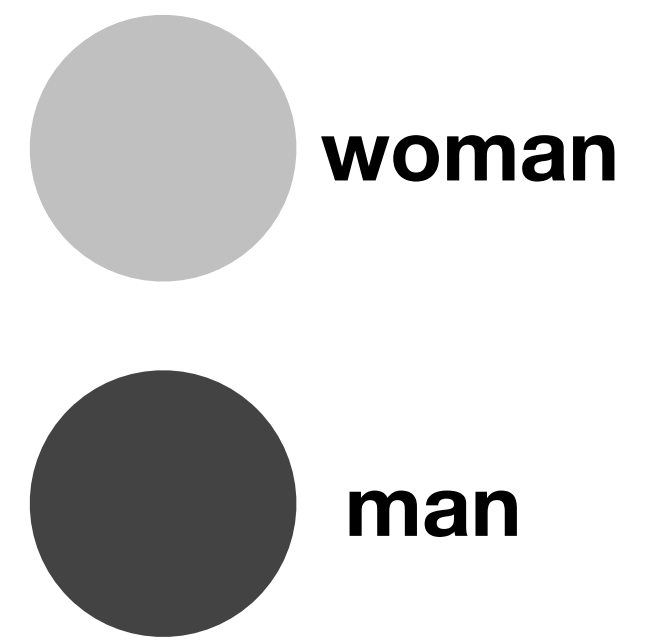




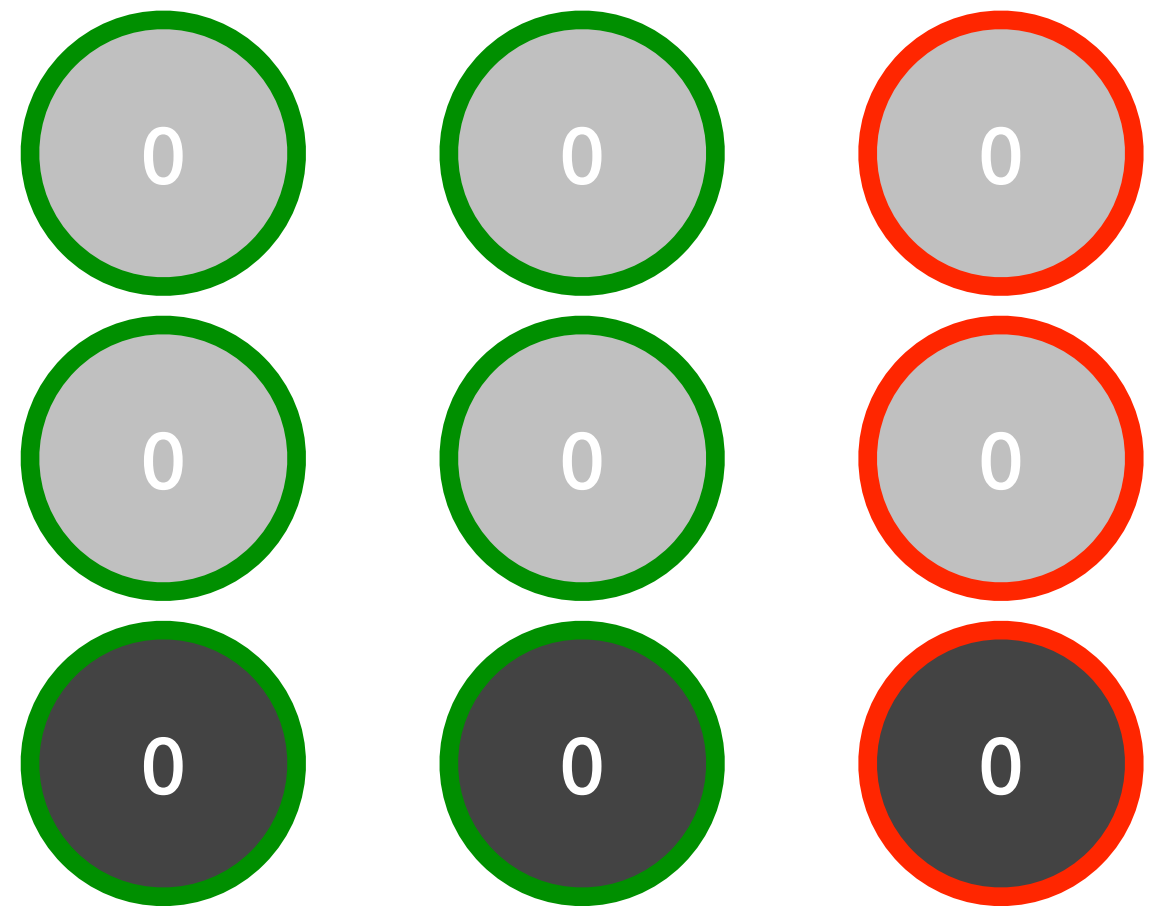
LOW SCORE



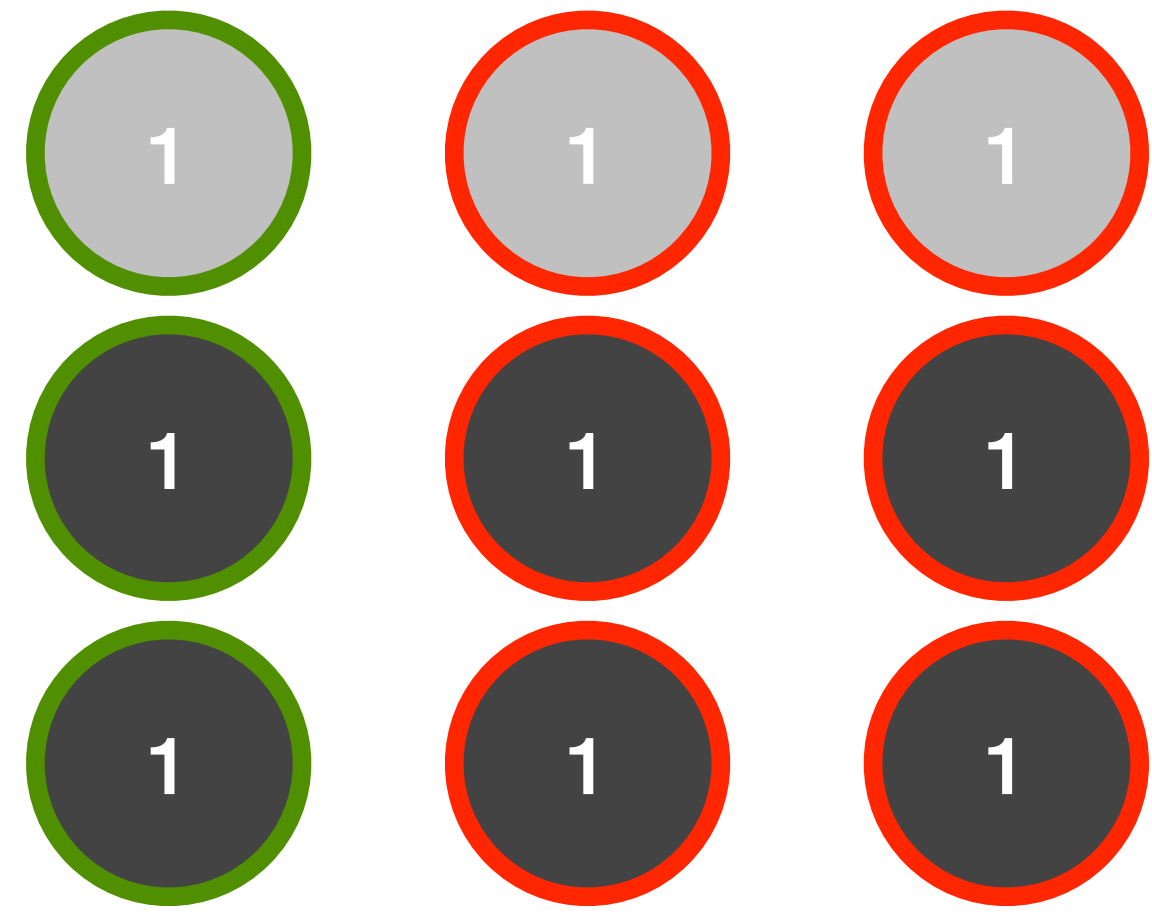
HIGH SCORE



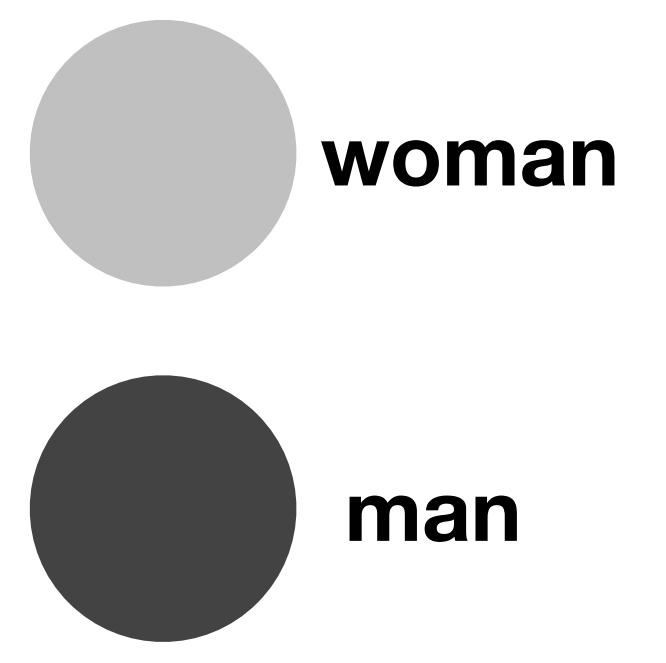
PREDICTION



LOW SCORE



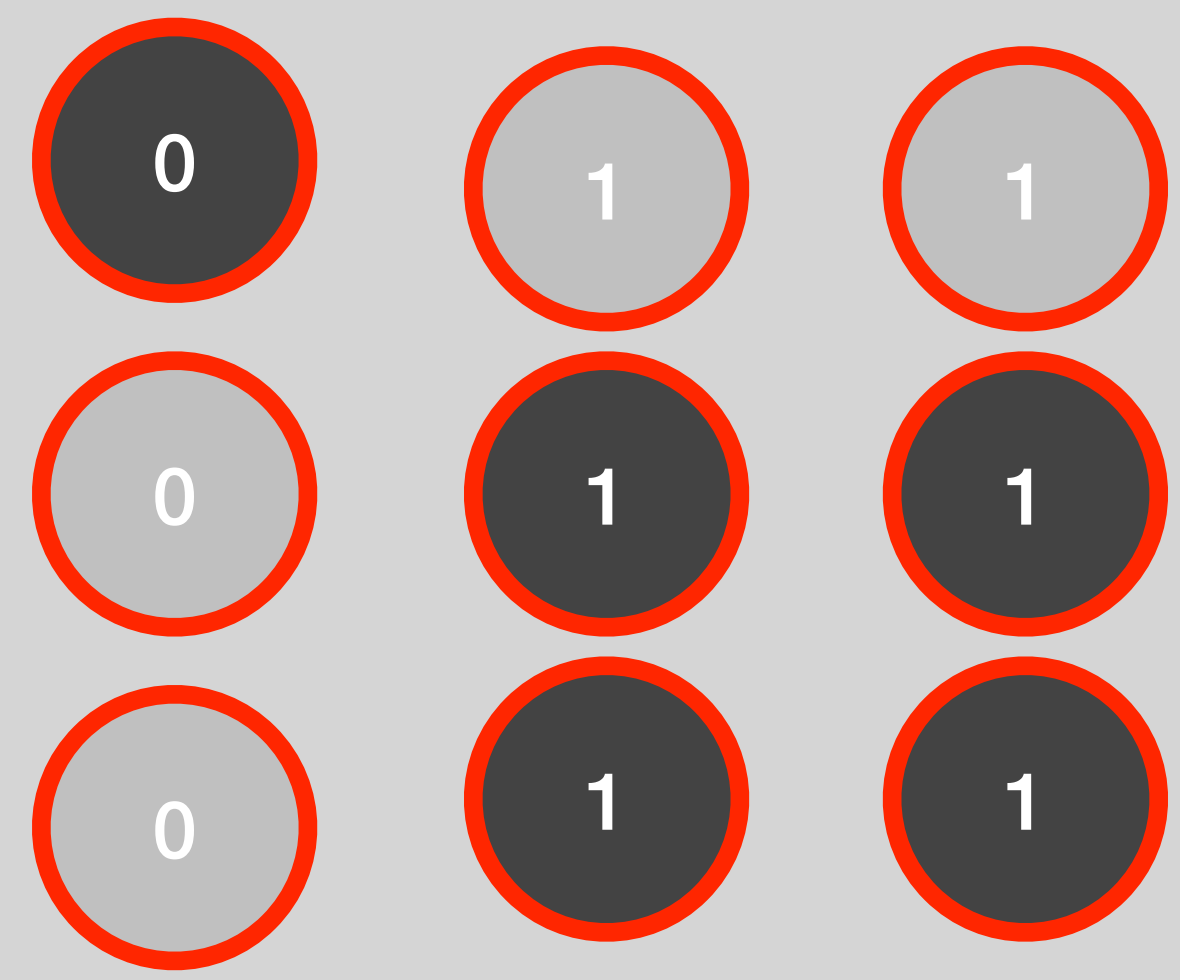
HIGH SCORE



POST-DICTION

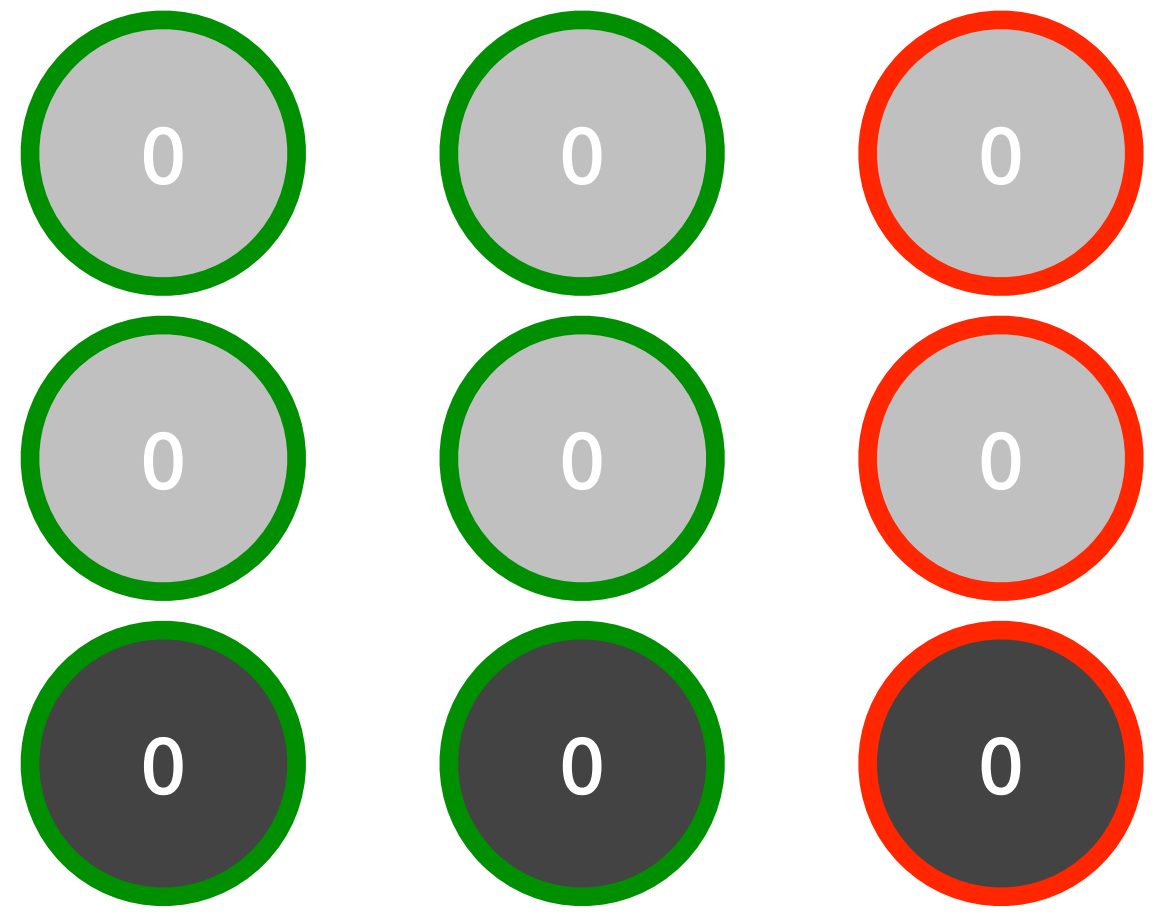


NO REOFFEND

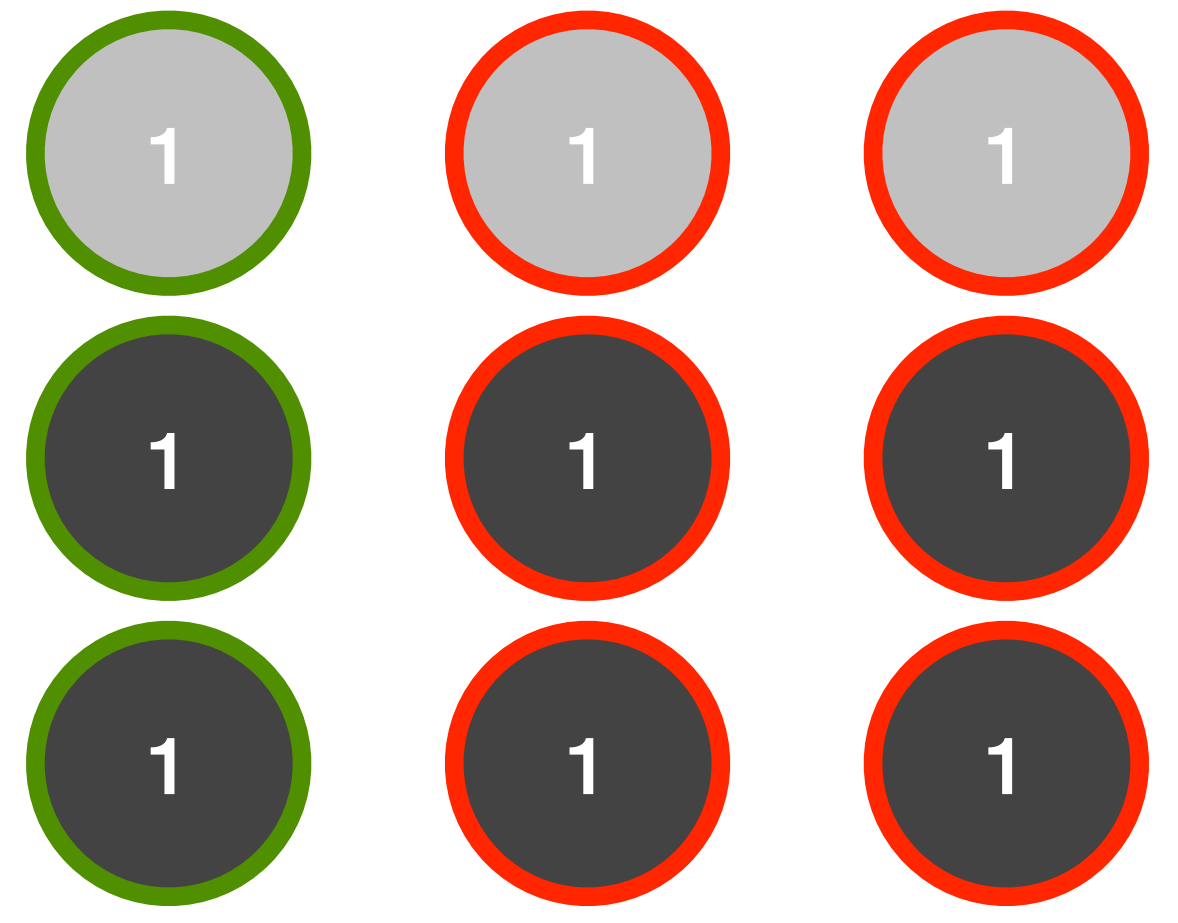


REOFFEND

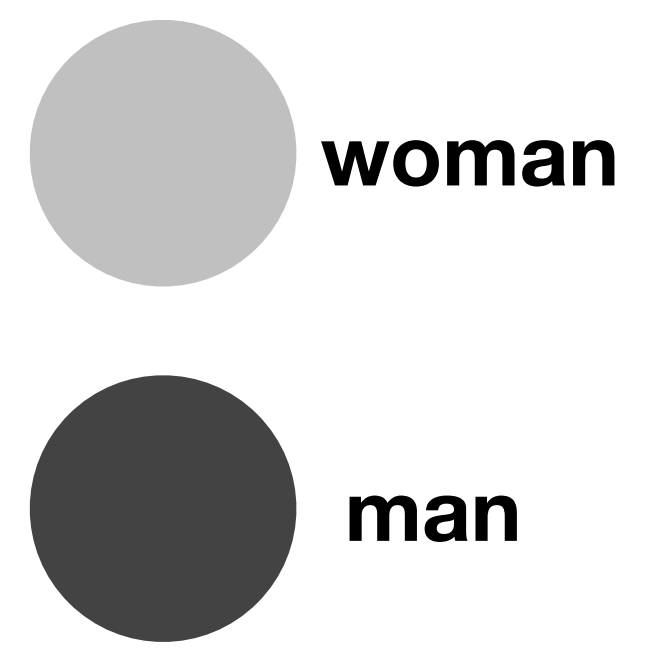
PREDICTION



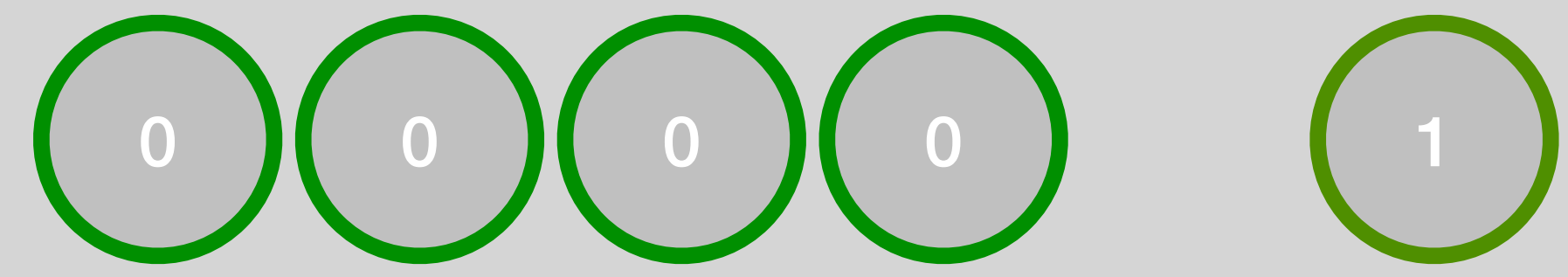
LOW SCORE



HIGH SCORE



POST-DICTION



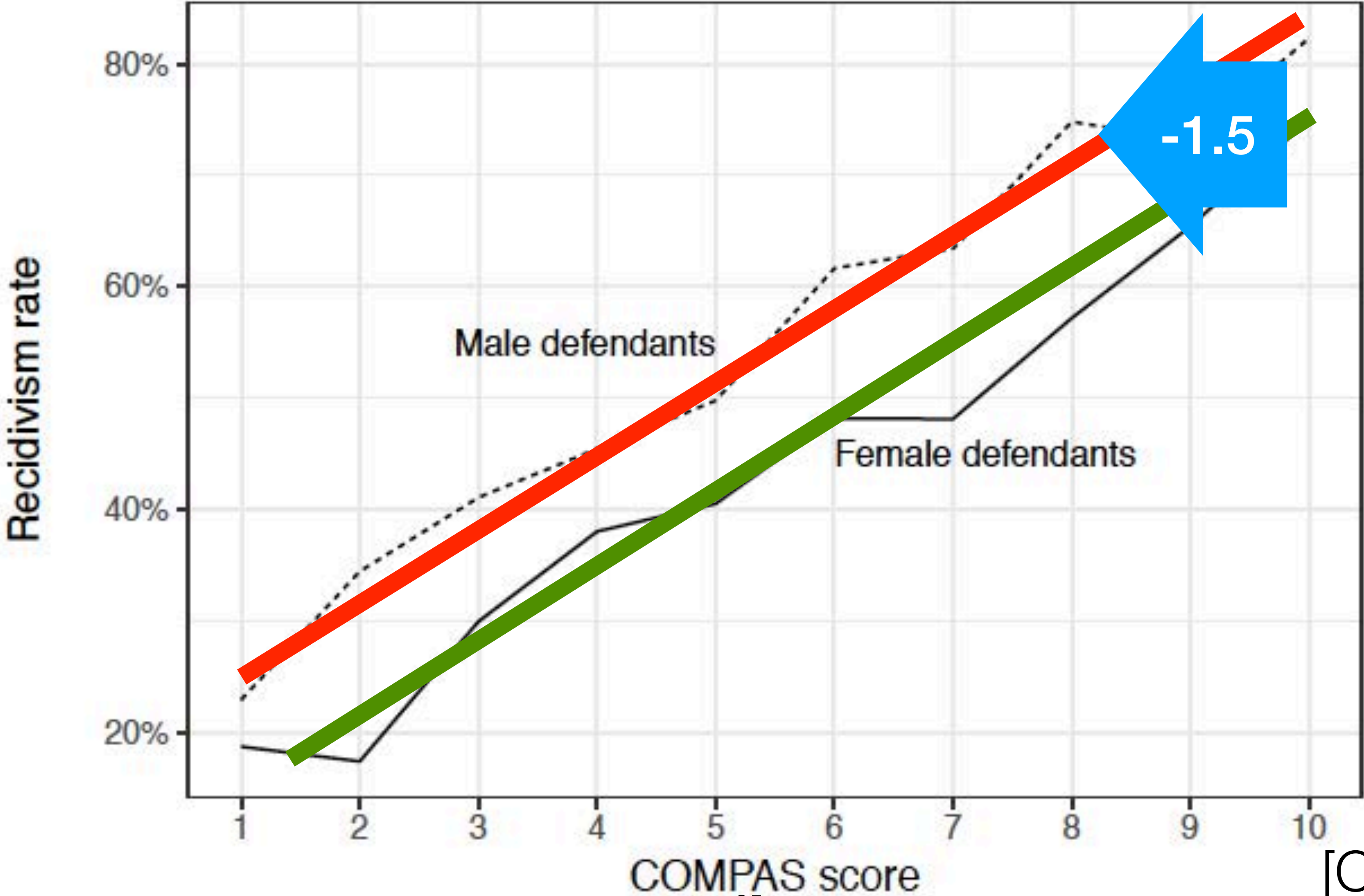
NO REOFFEND



REOFFEND



Male vs females...



[Corbett-Davies 2018]

Learnings

- Individual vs group
- Protected attributes
- Definitions of fairness
- Cannot be fair according to multiple definitions
- Society decides criterion, AI implements it

Face analysis

Face recognition

We Teach A.I. Systems Everything, Including Our Biases

Facial Recognition Is Accurate, if You're a White Guy

By Steve Lohr

Feb. 9, 2018



TECH

WORK LIFE

CREATIVITY

IMPACT

AUDIO

FAST COMPANY

C O M P A S S

05.22.19

Sign up for our daily email.

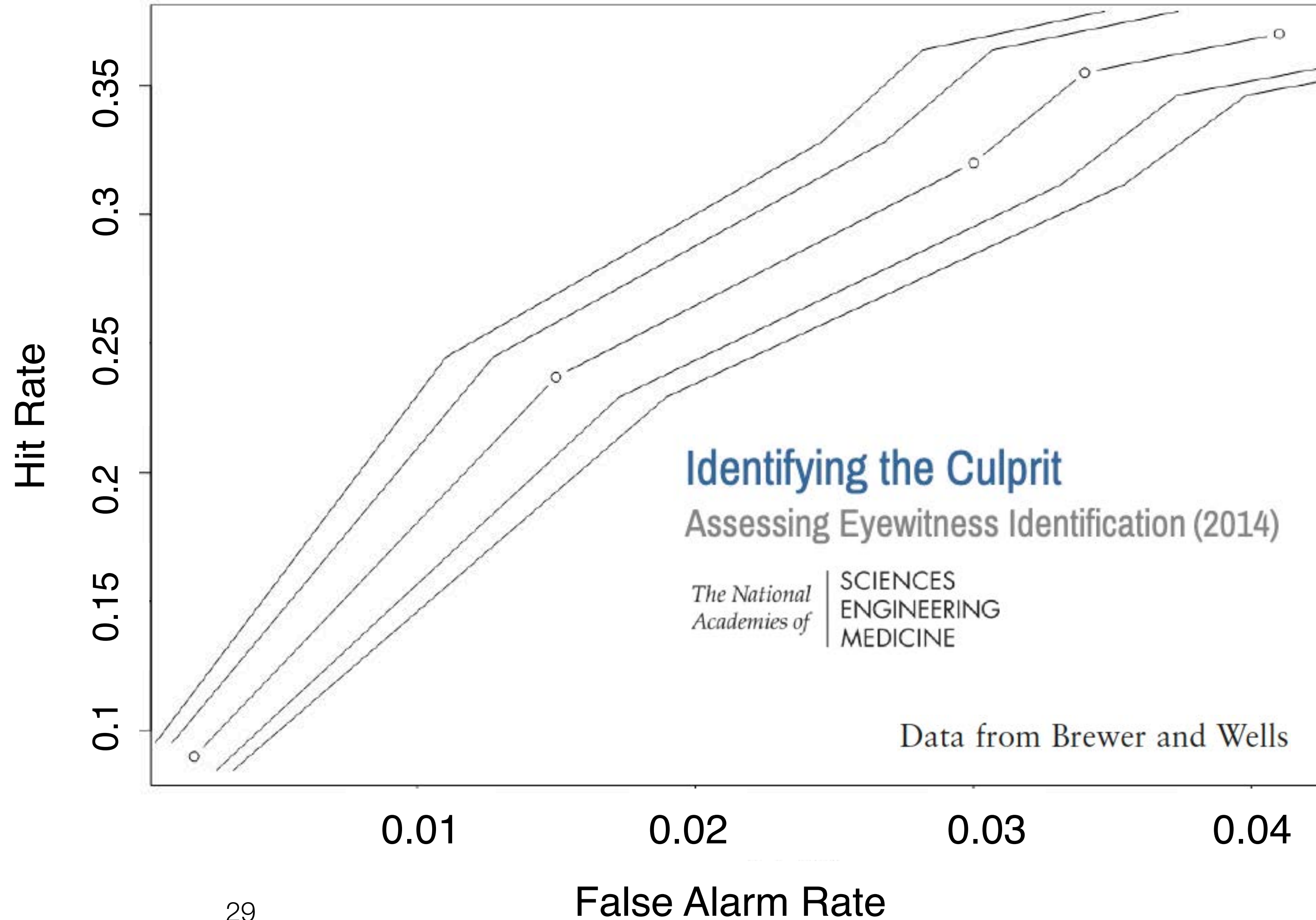
Enter your email

Here's AOC calling out the vicious circle of white men building biased face AI

Human error



Ronald Cotton and Jennifer Thompson

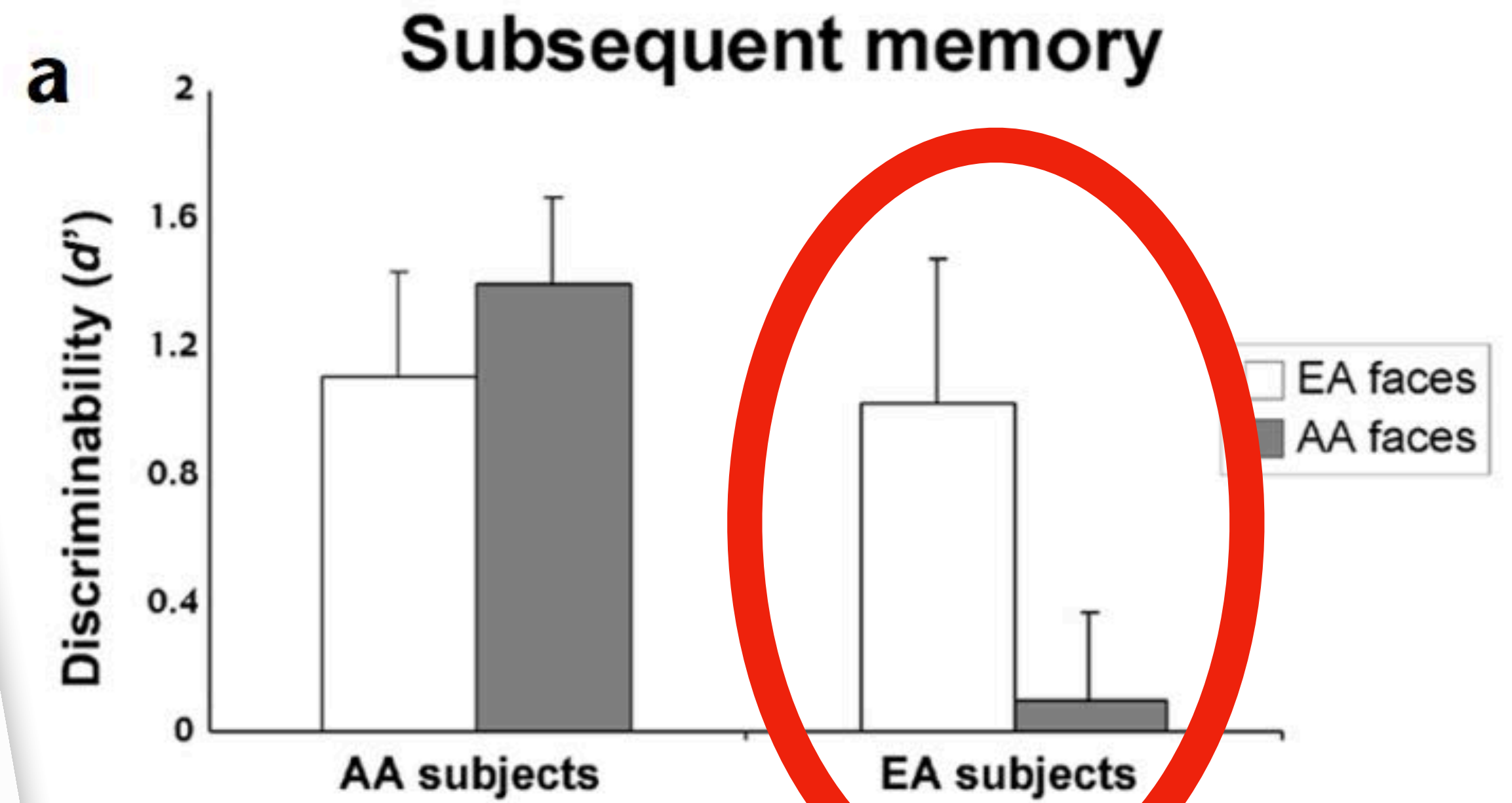


Human bias

Differential responses in the fusiform region to same-race and other-race faces

Alexandra J. Golby, John D. E. Gabrieli, Joan Y. Chiao and Jennifer L. Eberhardt

Departments of Radiology and Psychology, Jordan Hall-Building 420, Stanford University, Stanford, California 94305, USA
Correspondence should be addressed to J.L.E. (jle@psych.stanford.edu)
nature neuroscience • volume 4 no 8 • august 2001



ECONOMIC VIEW

Biased Algorithms Are Easier to Fix Than Biased People

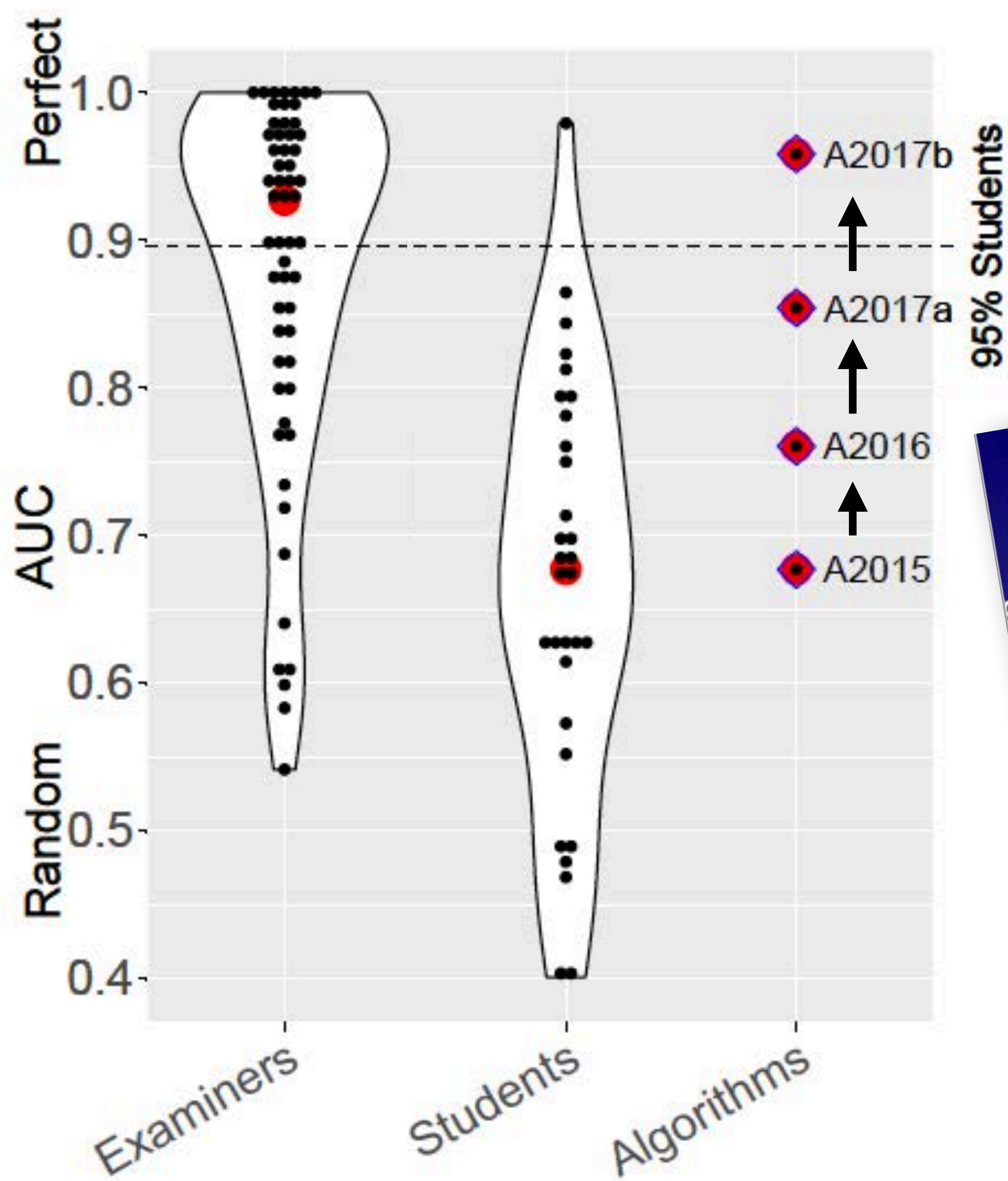
Racial discrimination by algorithms or by people is harmful — but that's where the similarities end.

By Sendhil Mullainathan

Dec. 6, 2019



Accuracy, consistency



PNAS

Face recognition accuracy of forensic examiners, superrecognizers, and face recognition algorithms

P. Jonathon Phillips^{a,1}, Amy N. Yates^a, Ying Hu^b, Carina A. Hahn^b, Eilidh Noyes^b, Kelsey Jackson^b, Jacqueline G. Cavazos^b, Géraldine Jeckeln^b, Rajeev Ranjan^c, Swami Sankaranarayanan^c, Jun-Cheng Chen^d, Carlos D. Castillo^d, Rama Chellappa^c, David White^e, and Alice J. O'Toole^b

^aInformation Access Division, National Institute of Standards and Technology, Gaithersburg, MD 20899; ^bSchool of Behavioral and Brain Sciences, The University of Texas at Dallas, Richardson, TX 75080; ^cDepartment of Electrical and Computer Engineering, University of Maryland Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20854; ^dUniversity of Maryland Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20854; and ^eSchool of Psychology, The University of New South Wales, Sydney, NSW 2052, Australia



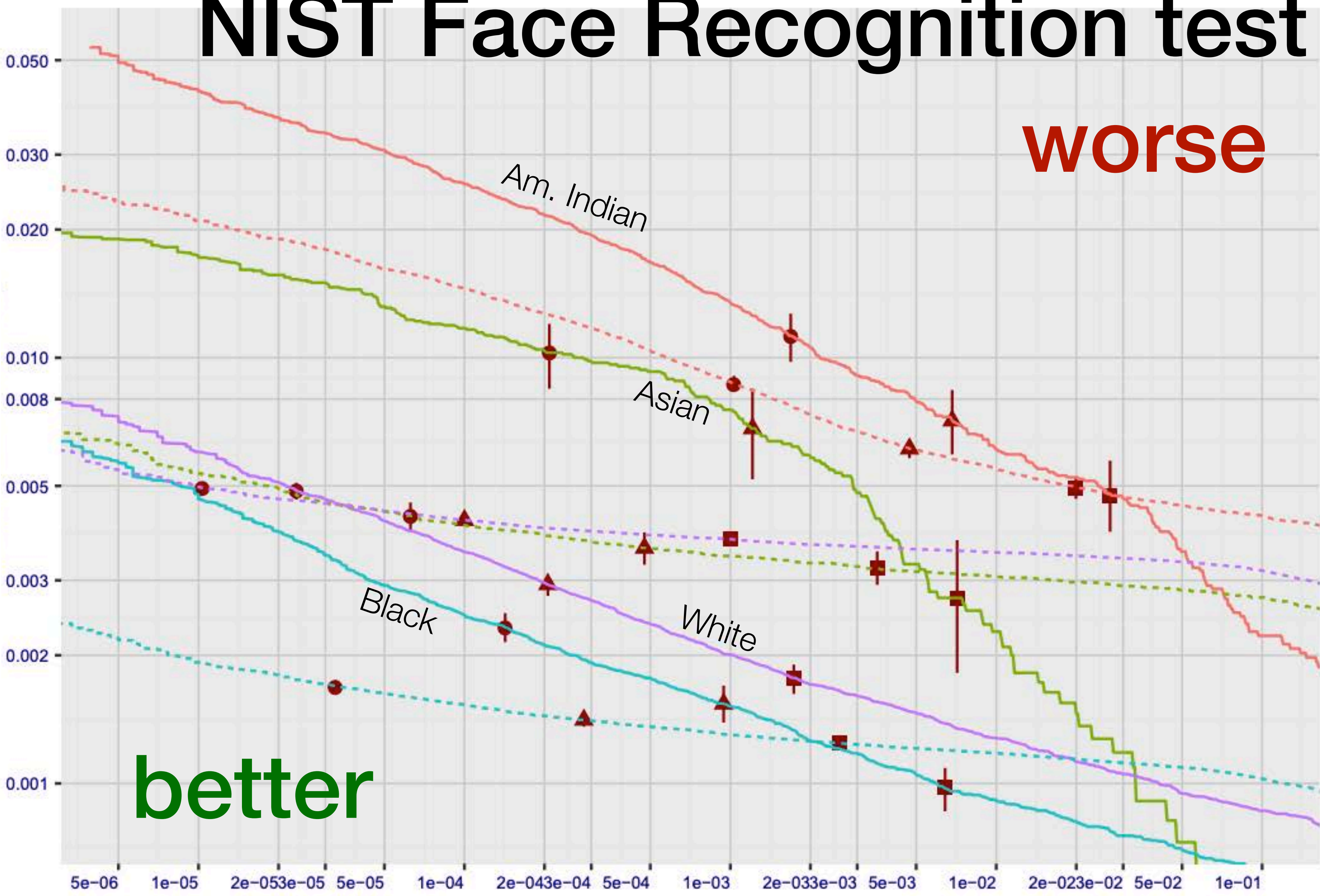
Benefits

- Cost
- Speed
- Scalability
- Fairness
- Accuracy
- Accountability
- Improvability
- Consistency

NIST Face Recognition test

(MUGSHOT)

False reject rate



False accept rate

[NIST 2019]

Race Bias ?

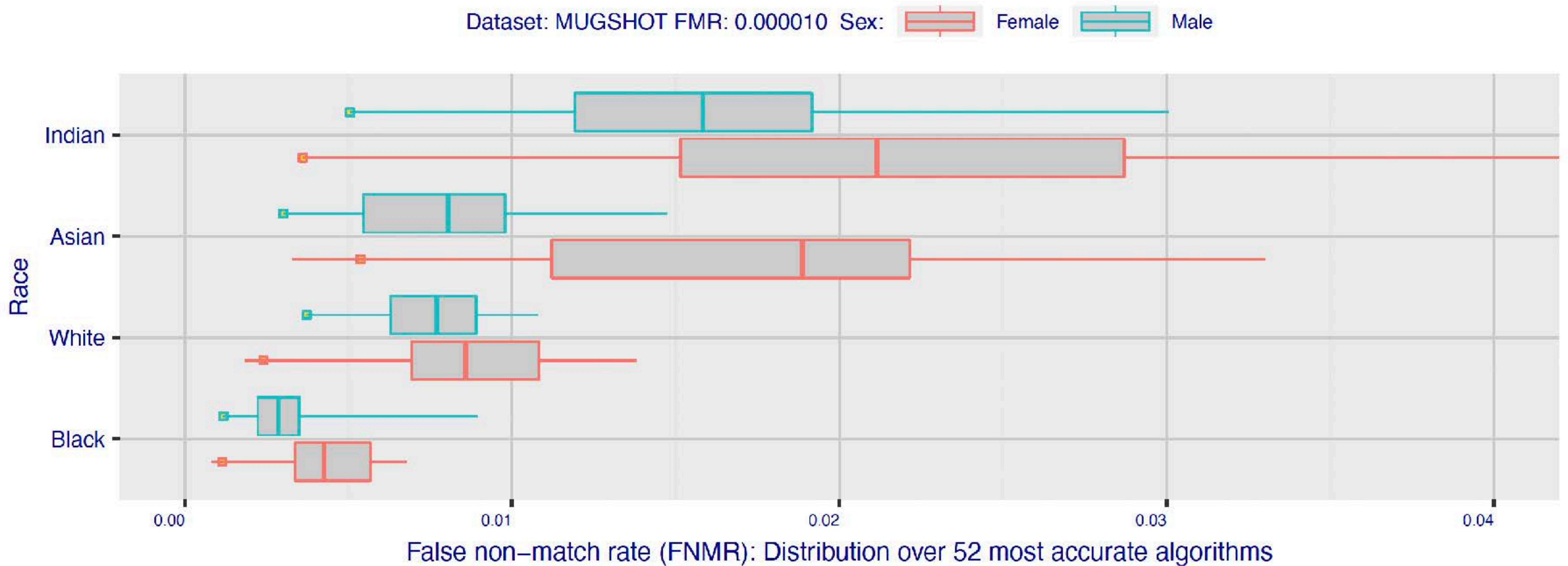
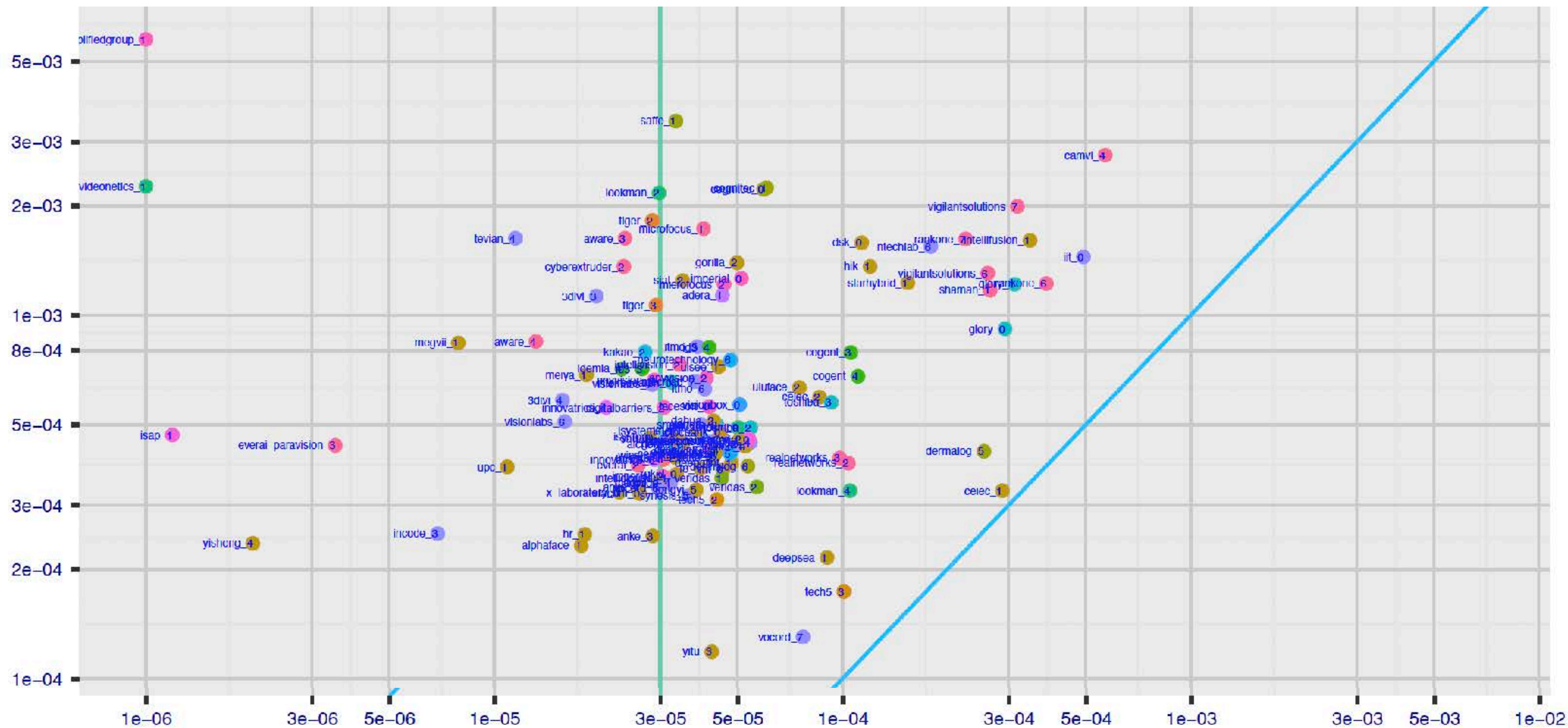


Figure 17: For mugshot comparisons, the figure shows the distribution of FNMR values over the 52 most accurate verification algorithms, by sex and race. The threshold was set for each algorithm to achieve $FMR = 0.00001$ over all imposter comparisons. The line within each box is the median over those algorithms; the box itself spans the interquartile range (26 algorithms) and the lines here extend to minimum and maximum values. The small box on the left side indicates the accuracy for best algorithm overall, on this dataset `alphaface3601`.

Race bias ?

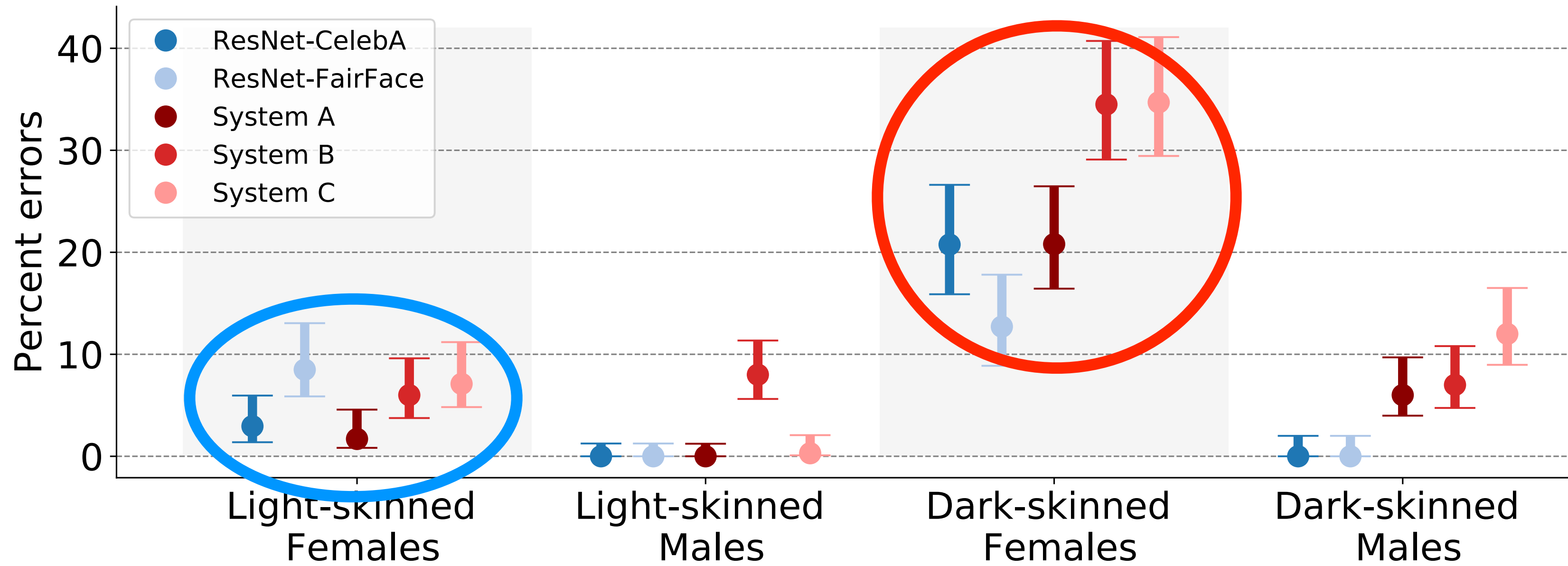
False accept rate in men from W. Africa

W. Africa



False accept rate in men from E. Europe

Gender classification



PPB
dataset
(Scandinavia +
Sub-Saharan Africa)
[Buolamwini & Gebru 2018]



Parliament web sites of
Sweden, Finland
and Iceland

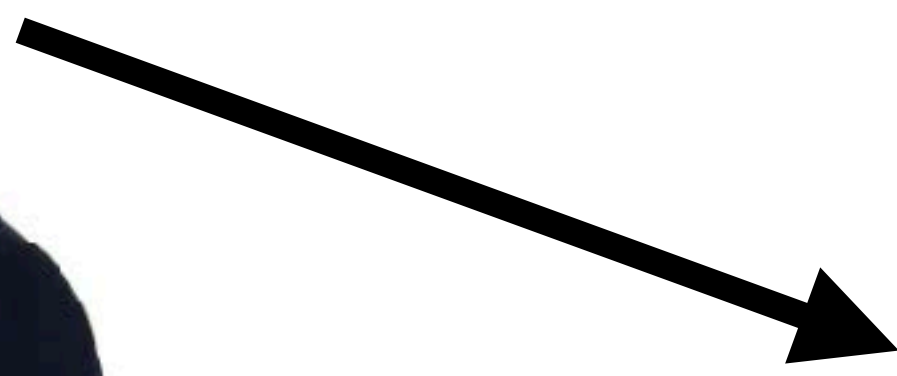
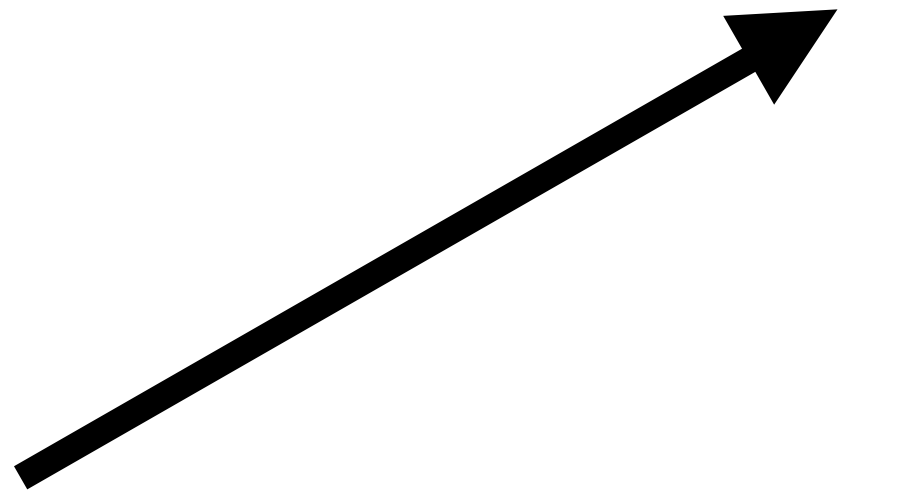
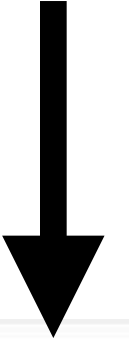


Parliament web sites of
Rwanda, Senegal
and South Africa

Observational study



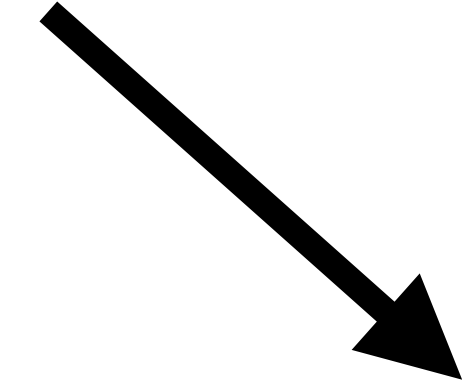
sample



annotate

algorithm

Female
Dark Skin



prediction

	Light skin	Dark skin
Female	3%	25%
Male	5%	3%

The New York Times

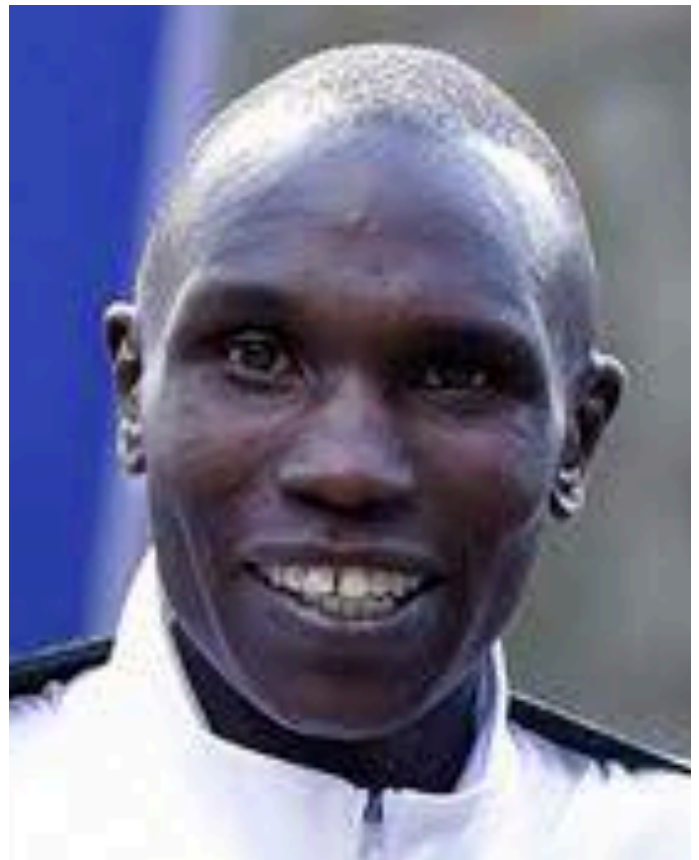
Facial Recognition Is Accurate, if You're a White Guy

By Steve Lohr

Feb. 9, 2018



COVID mortality

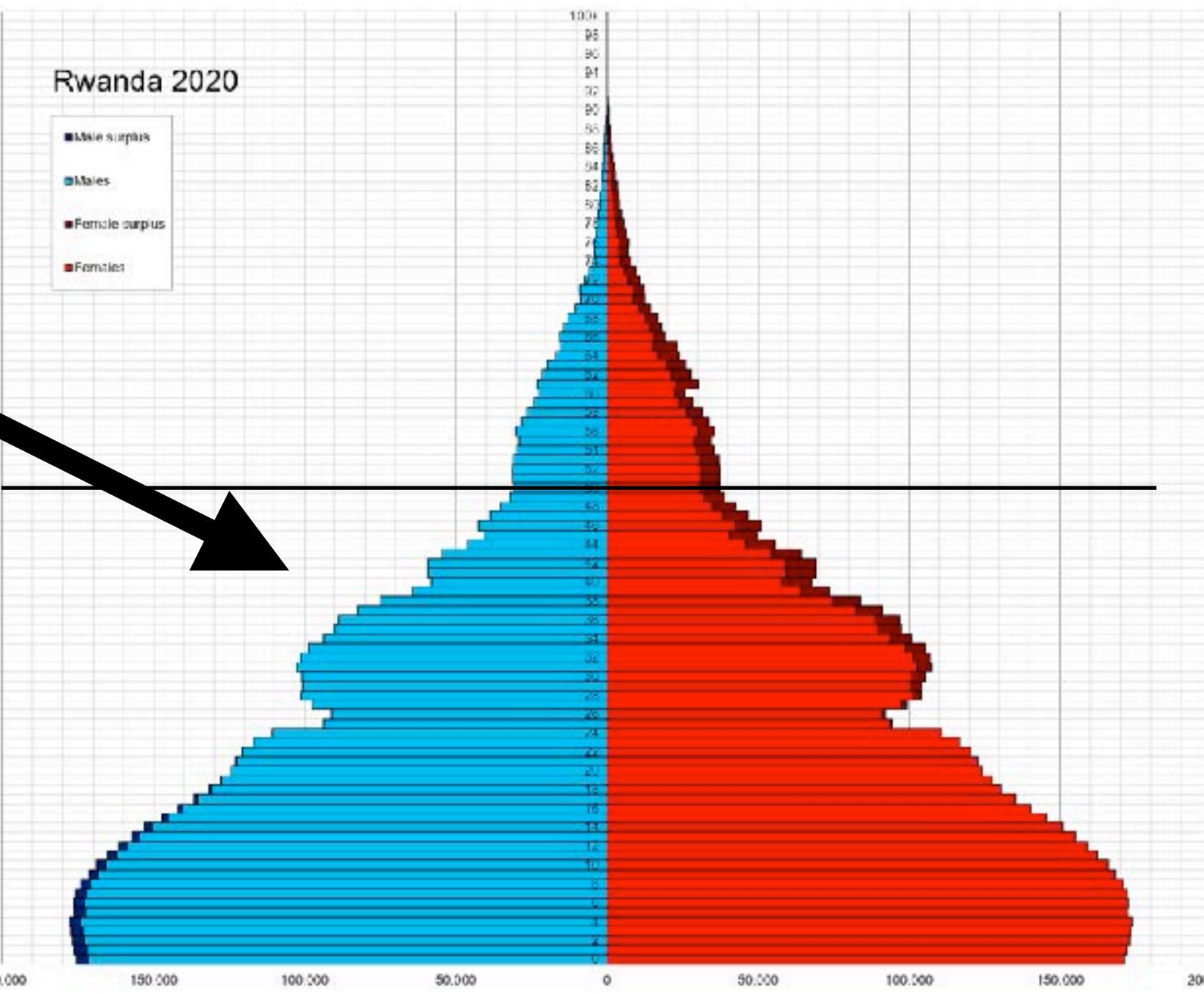
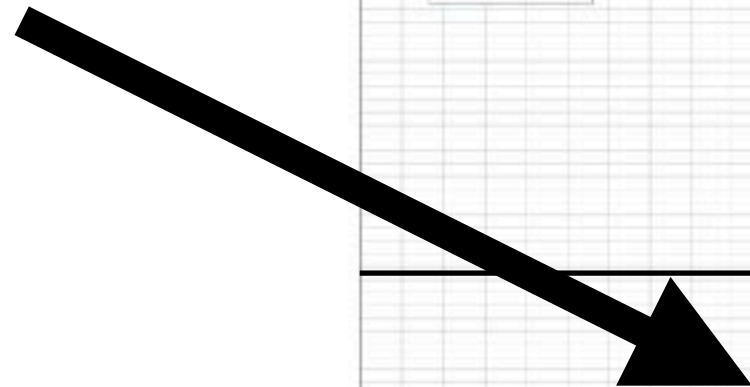
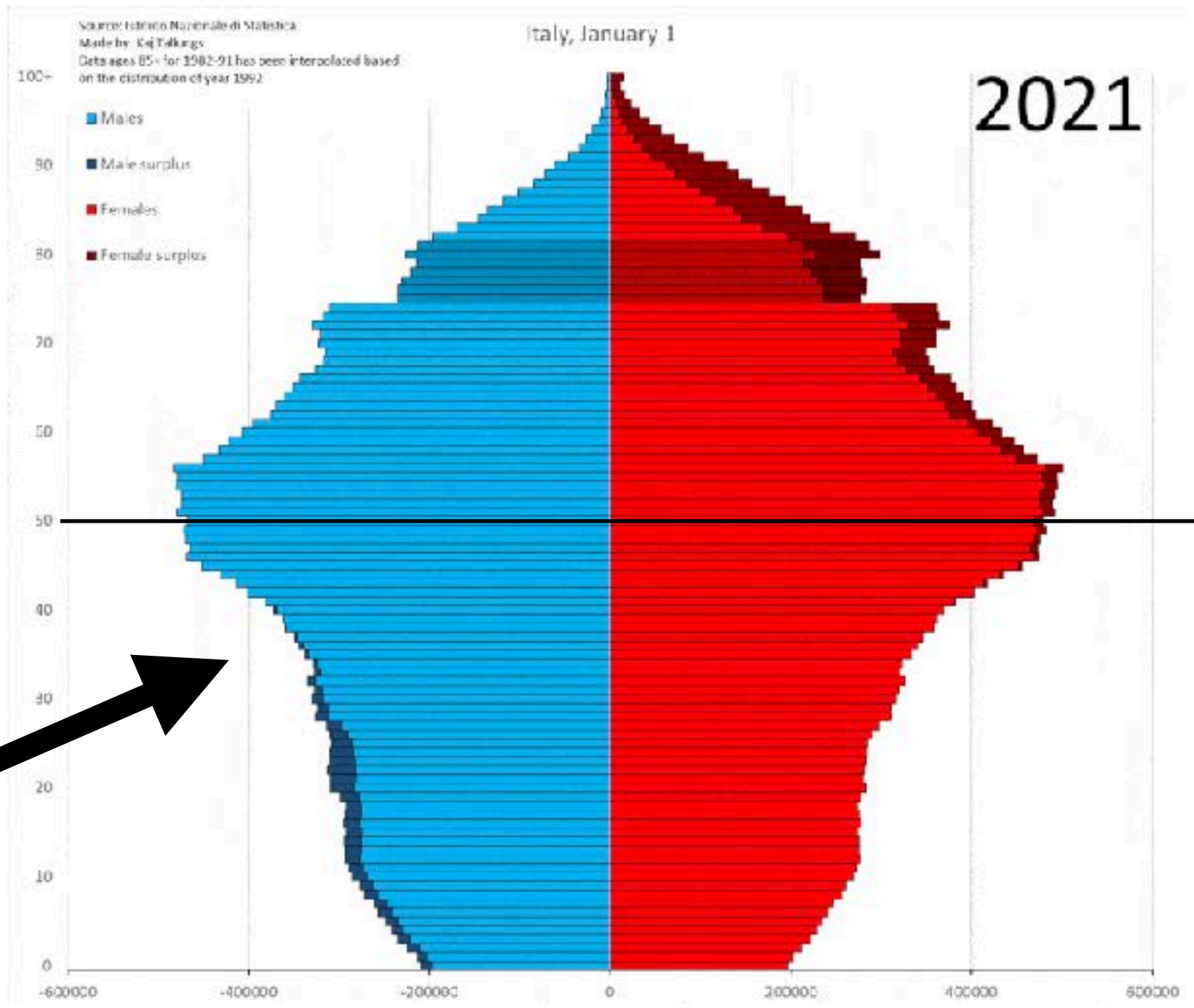
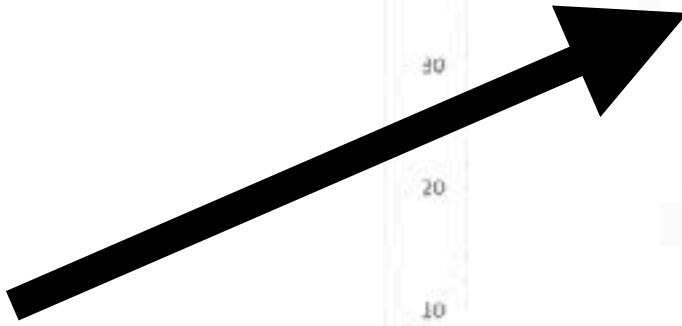


COVID MORTALITY

Country	Deaths/100K pop.
United Kingdom	
Poland	
Italy	
Sudan	
Rwanda	
Kenya	

[Johns Hopkins Coronavirus Resource Center, 2023]

Huge bias vis-a-vis ethnicity?



The NEW ENGLAND JOURNAL of MEDICINE
**Chocolate Consumption, Cognitive Function,
 and Nobel Laureates**
 Franz H. Messerli, M.D.

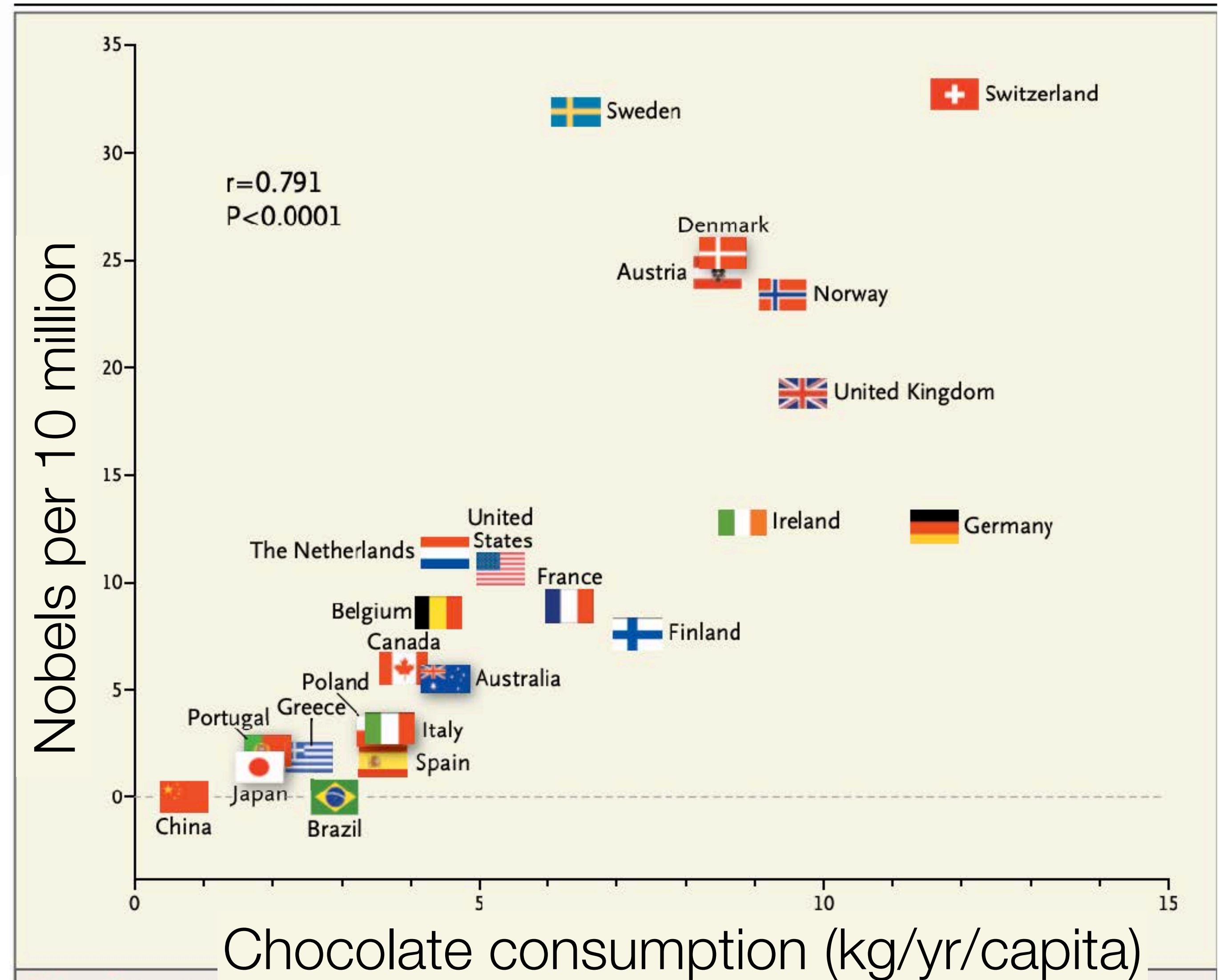
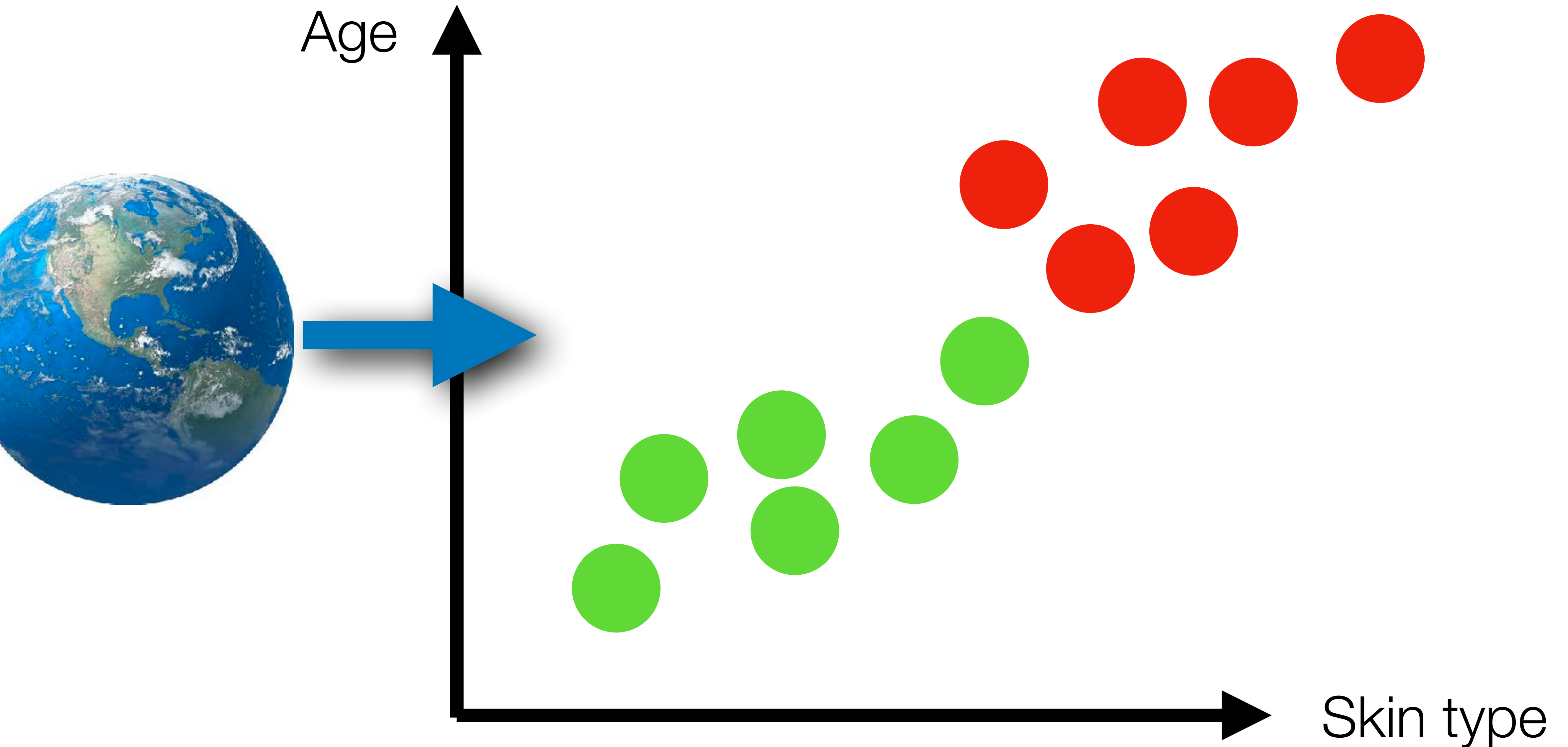
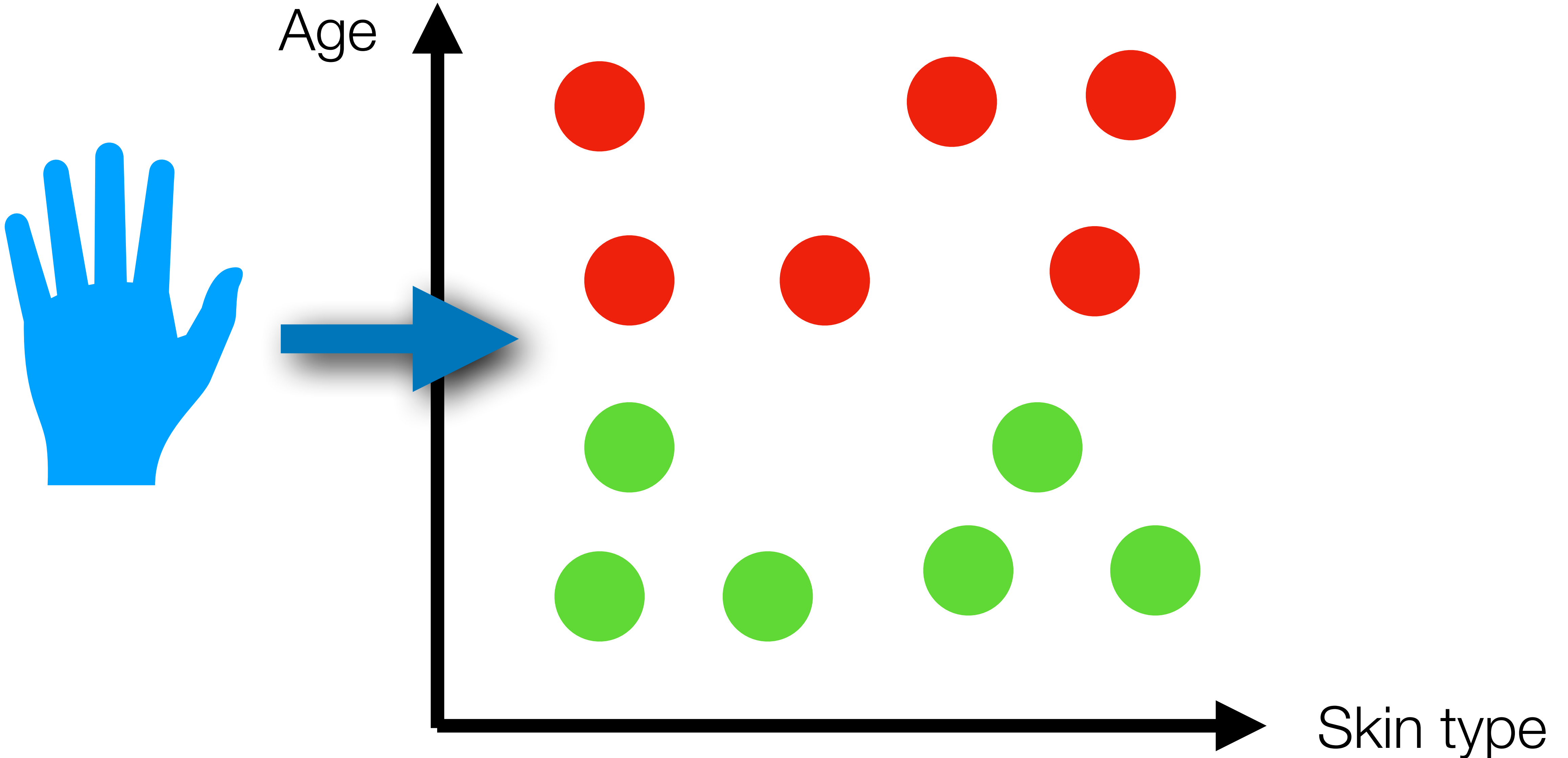


Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

Observational study



Experiment



ARE EMILY AND GREG MORE EMPLOYABLE
THAN LAKISHA AND JAMAL?
A FIELD EXPERIMENT ON LABOR MARKET DISCRIMINATION

Marianne Bertrand
Sendhil Mullainathan

Working Paper 9873
<http://www.nber.org/papers/w9873>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
July 2003

ARE EMILY AND GREG MORE EMPLOYABLE
 THAN LAKISHA AND JAMAL?
 A FIELD EXPERIMENT ON LABOR MARKET DISCRIMINATION

Marianne Bertrand
 Sendhil Mullainathan

Working Paper 9873
<http://www.nber.org/papers/w9873>

NATIONAL BUREAU OF ECONOMIC RESEARCH
 1050 Massachusetts Avenue
 Cambridge, MA 02138
 July 2003

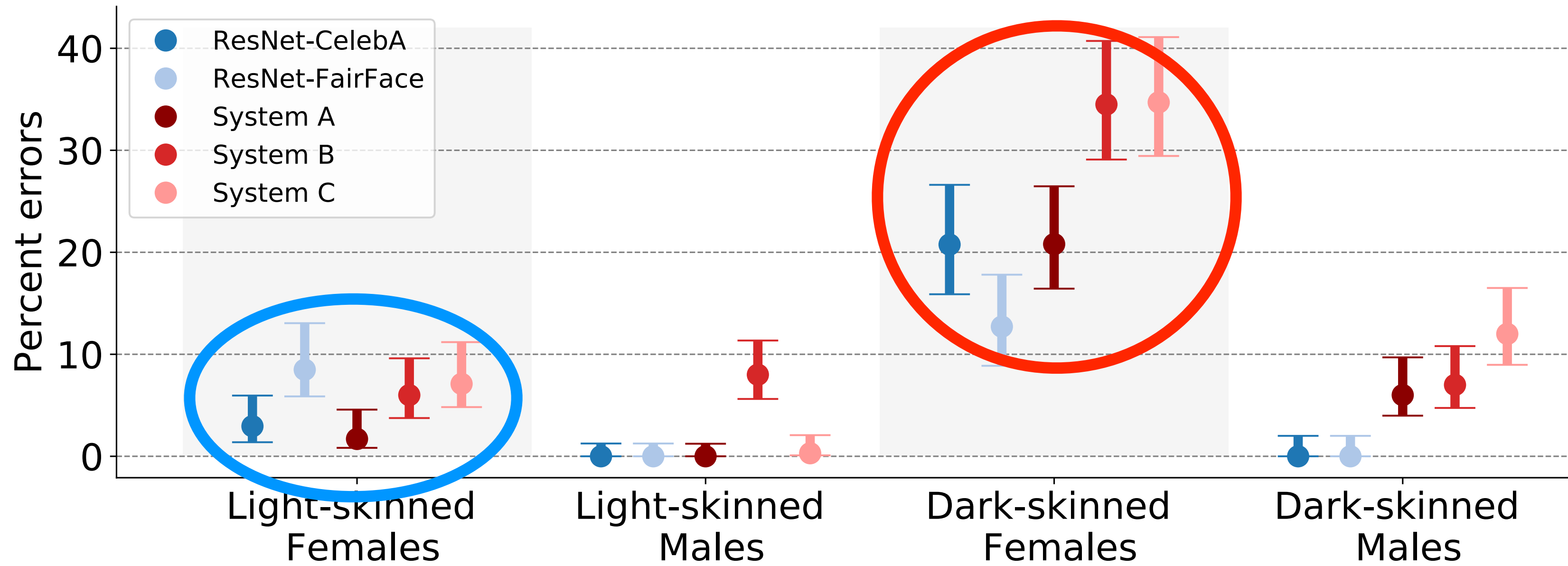
Table 1
Mean Call-Back Rates By Racial Soundingness of Names ^a

	<i>Call-Back Rate for White Names</i>	<i>Call-Back Rate for African American Names</i>	<i>Ratio</i>	<i>Difference (p-value)</i>
Sample:				
All sent resumes	10.06% [2445]	6.70% [2445]	1.50	3.35% (.0000)
Chicago	8.61% [1359]	5.81% [1359]	1.48	2.80% (.0024)
Boston	11.88% [1086]	7.83% [1086]	1.52	4.05% (.0008)
Females	10.33% [1868]	6.87% [1893]	1.50	3.46% (.0001)
Females in administrative jobs	10.93% [1363]	6.81% [1364]	1.60	4.12% (.0001)
Females in sales jobs	8.71% [505]	6.99% [529]	1.25	1.72% (.1520)
Males	9.19% [577]	6.16% [552]	1.49	3.03% (.0283)

Observation -> correlation

Experiment -> cause

Gender classification



PPB
dataset
(Scandinavia +
Sub-Saharan Africa)
[Buolamwini & Gebru 2018]

Gender-typical



Non-typical



[Mutukhumar et al. 2018]

Experimental method in computer vision?

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

ECCV2020

Benchmarking Algorithmic Bias in Face Recognition:
An Experimental Approach Using Synthetic Faces and Human Evaluation

Hao Liang
Rice University
h1106@rice.edu

Pietro Perona
California Institute of Technology and AWS
perona@caltech.edu, peronapp@amazon.com

Guha Balakrishnan
Rice University
guha@rice.edu

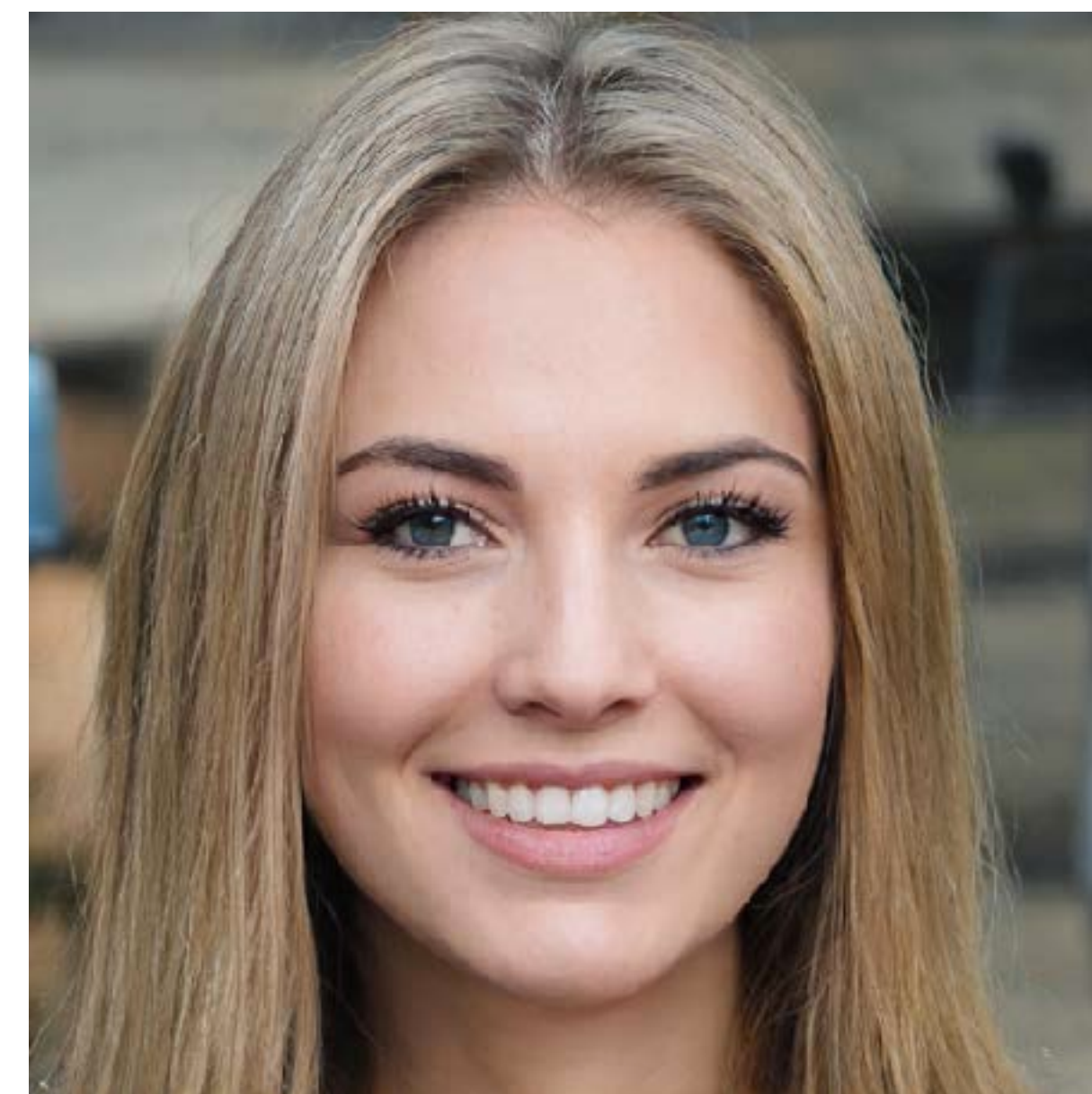
ICCV2023



Darker skin

Idea 1: use synthetic faces

<https://www.thispersondoesnotexist.com/>



Idea 2: Control of latent space

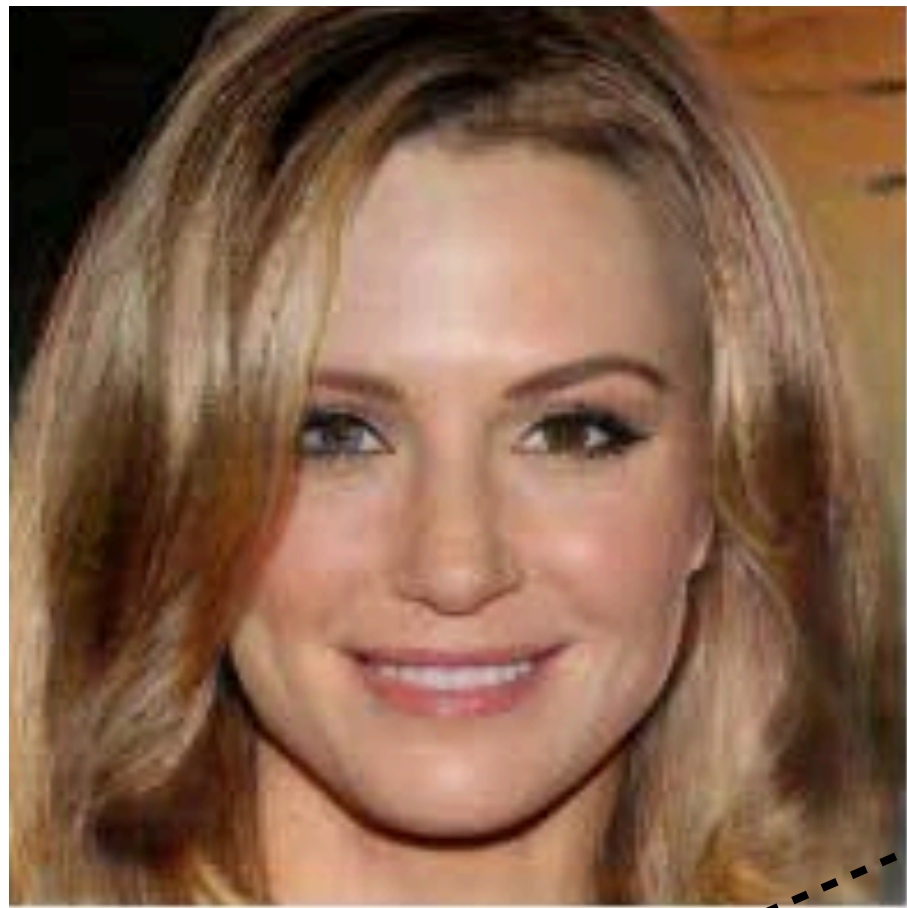
Interpreting the Latent Space of GANs for Semantic Face Editing

Yujun Shen¹, Jinjin Gu², Xiaoou Tang¹, Bolei Zhou¹

¹The Chinese University of Hong Kong ²The Chinese University of Hong Kong, Shenzhen
{sy116, xtang, bzhou}@ie.cuhk.edu.hk, jinjingu@link.cuhk.edu.cn

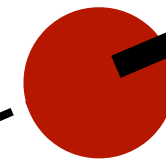
arXiv:1907.10786v3 [cs.CV] 31 Mar 2020

female



Z_1

female

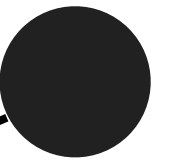


male



male

Z_2



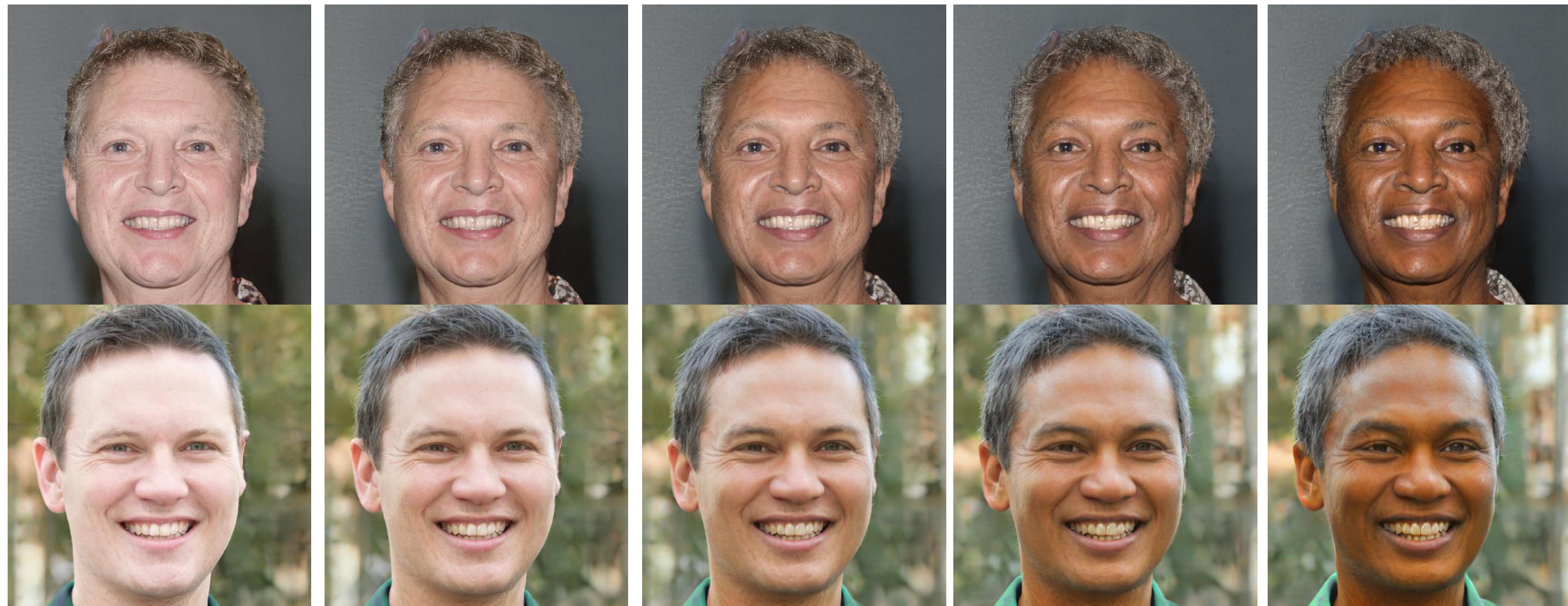
← Age →

← Facial Hair →



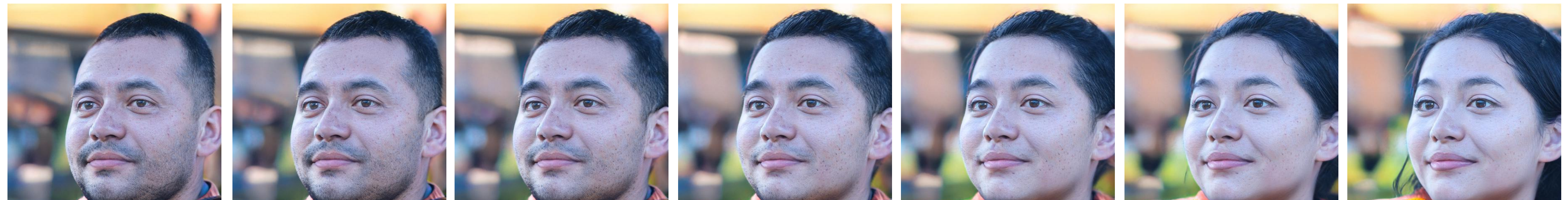
← Skin Color →

← Hair Length →



← Hair Length →

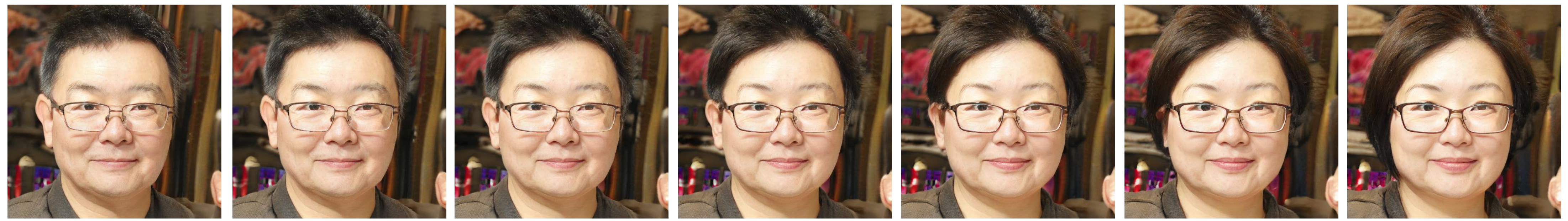
No Orth.



Orth.



No Orth.



Orth.





Skin Color

No Orth.



Orth.



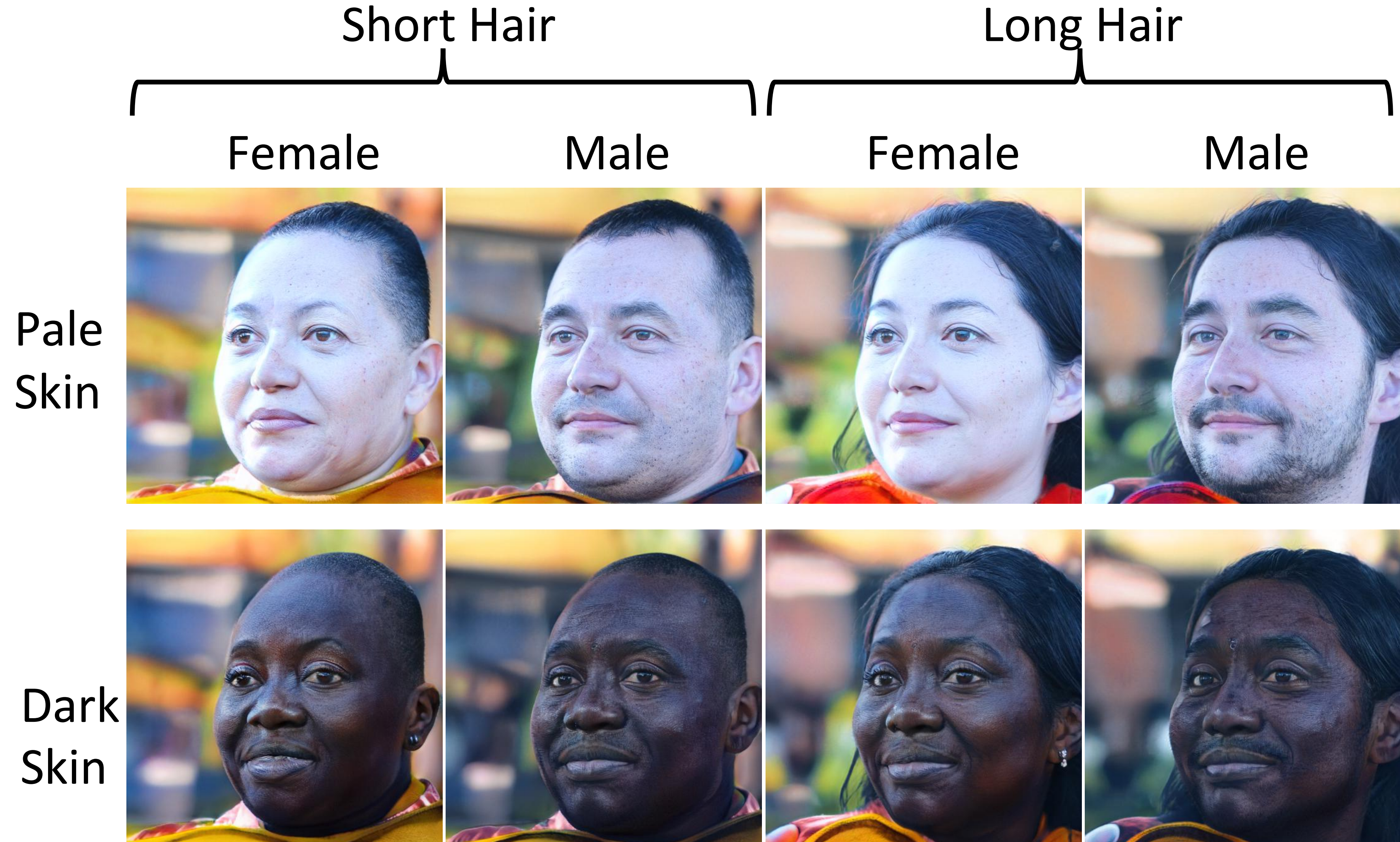
No Orth.



Orth.



Idea 3: matched samples

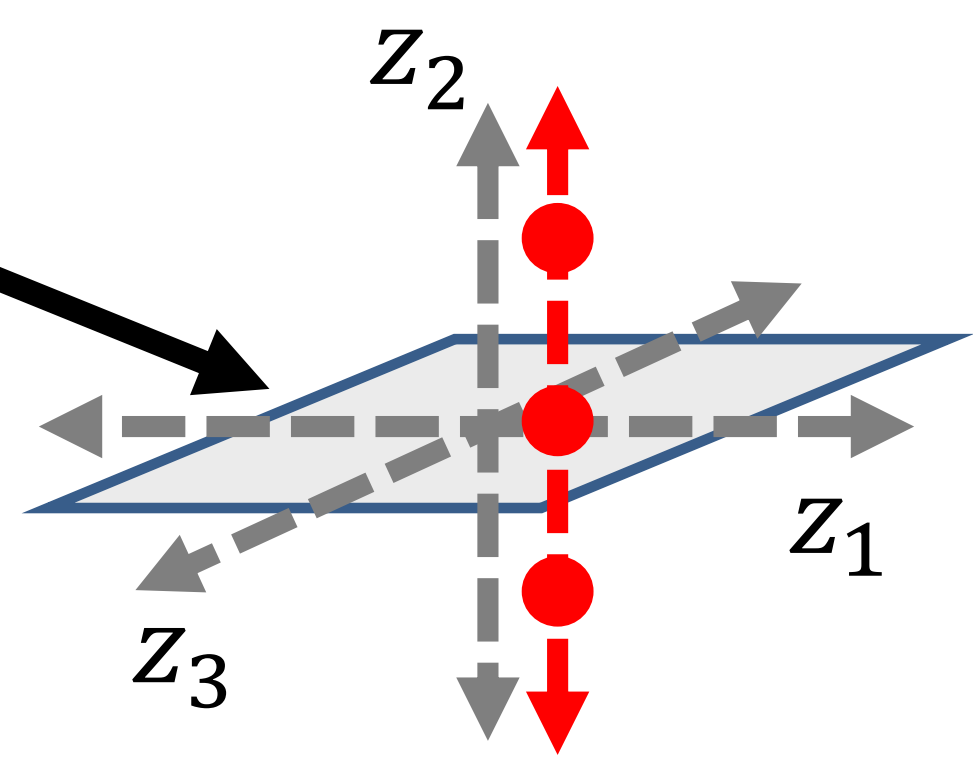


**Idea 4: calibrate with human
annotators**

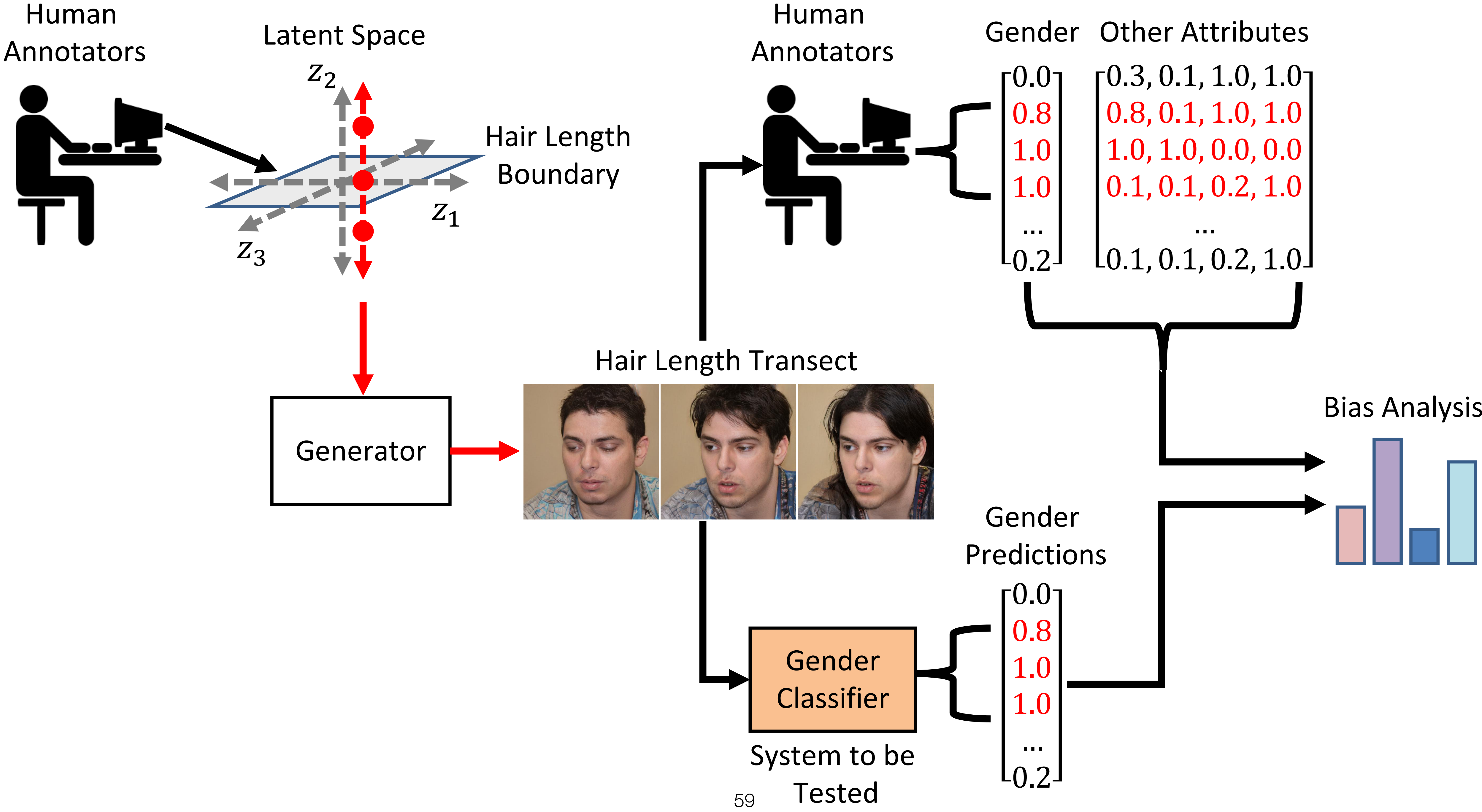
Human
Annotators



Latent Space



Hair Length
Boundary



Collecting and assessing human annotations

Instructions



[View full instructions](#)

[View tool guide](#)



I
Light, pale white



IV
Olive, mid brown



II
White, fair



V
Brown, dark brown



III
Medium white to olive



VI
Very dark brown, black

NOTE: Please do not worry if the face looks a little weird!

What's the skin color of this person? (Please refer to the scale on the left.)



Select an option

Light	1
Fair	2
Medium	3
Olive	4
Brown	5
Black	6



**SageMaker
Ground Truth**

Easily label training data for machine learning at scale

Instructions



[View full instructions](#)

[View tool guide](#)

* **Very short:** Bald or buzz cut. Between 0 and 1" (0-2.5cm)

* **Short:** 1-2" (2.5-5cm). Not easy to comb.

* **Medium:** May be combed, ears partially covered. 2-4" (5-10cm).

* **Long:** Covers ears, may reach shoulders. 4-8" (10-20cm).

* **Very long.** Beyond shoulder length.

What's the length of the hair of the person in this photograph?



Select an option

Very short 1

Short 2

Medium 3

Long 4

Very long 5



**SageMaker
Ground Truth**

Easily label training
data for machine
learning at scale

Instructions



[View full instructions](#)

[View tool guide](#)

* **Female:** Everyone will agree.

* **Probably female:** Most will say female. A few may say male.

***In between:** Could be male or female or gender-neutral.

***Probably male:** Most will say male. A few might say female.

***Male.** Everyone will agree.

NOTE: Please do not worry if the face looks a bit weird!

What would casual observers think about the gender of this face?



Select an option

Female 1

Probably female 2

In between 3

Probably male 4

Male 5



Zoom in



Zoom out



Move



Fit image

63

Submit

Instructions



[View full instructions](#)

[View tool guide](#)

* **Child:** Between 0 and 12.
Before puberty.

* **Teen:** 13-18, in middle/high
school.

* **Young adult:** 19-30.
Attending college, first job.

* **Adult:** 30-50. Parent, mid-
career.

* **Middle age:** 50-65 - Late
career, kids in college,
approaching retirement.

* **Senior.** 65+. Often retired,
may be elderly.

NOTE: Please do not worry if
the face looks a bit weird.

What's the approximate age of this person?



Select an option

Child	1
Teen	2
Young adult	3
Adult	4
Middle age	5
Senior	6



Zoom in



Zoom out



Move



Fit image

64

Submit

Instructions



[View full instructions](#)

[View tool guide](#)

* **Real for sure:** Most people would agree that this is real.

* **Likely real:** Many people would think it's real, but some would have doubts.

* **Uncertain:** People could come out either way.

* **Likely fake:** Many would say it's fake, but not everyone.

* **Fake:** Most people would not think that this is real.

Is this the photograph of a real face, or a fake generated by a computer?



Select an option

Real for sure 1

Likely real 2

Uncertain 3

Likely fake 4

Fake for sure 5



Zoom in



Zoom out



Move



Fit image

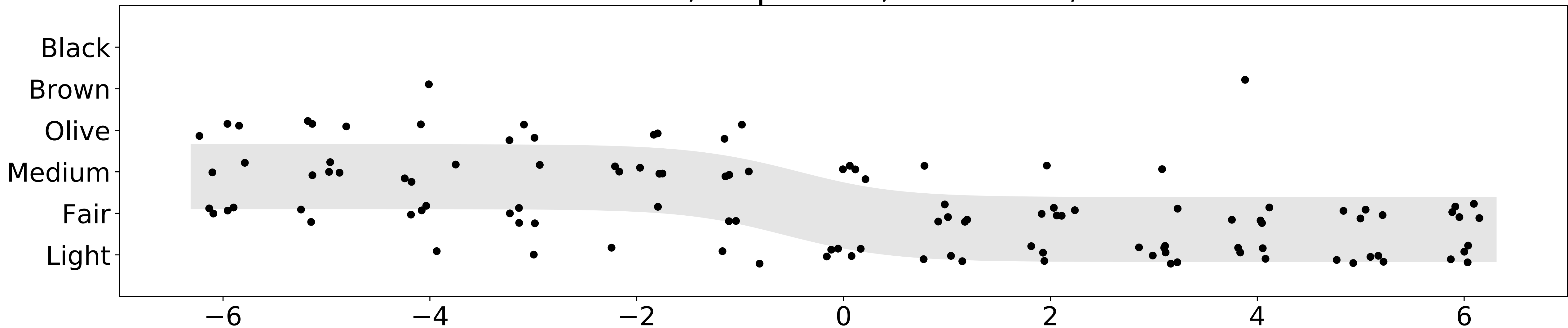
65

Submit

Skin tone



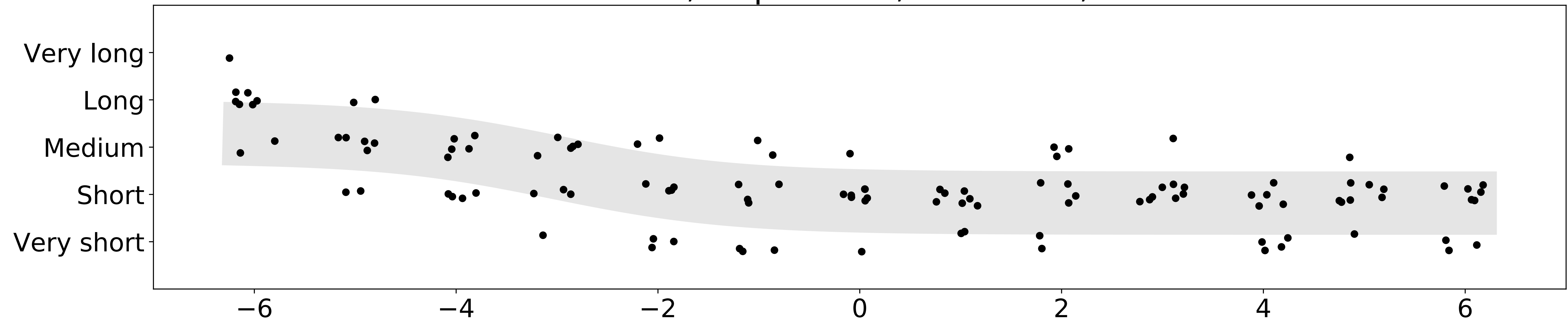
fit: shift=-0.51, slope=0.5, min=0.61, max=1.9



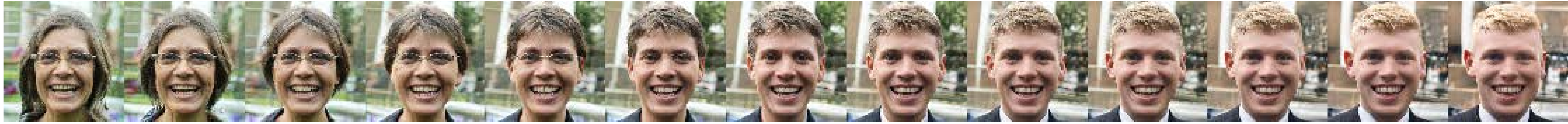
Hair length



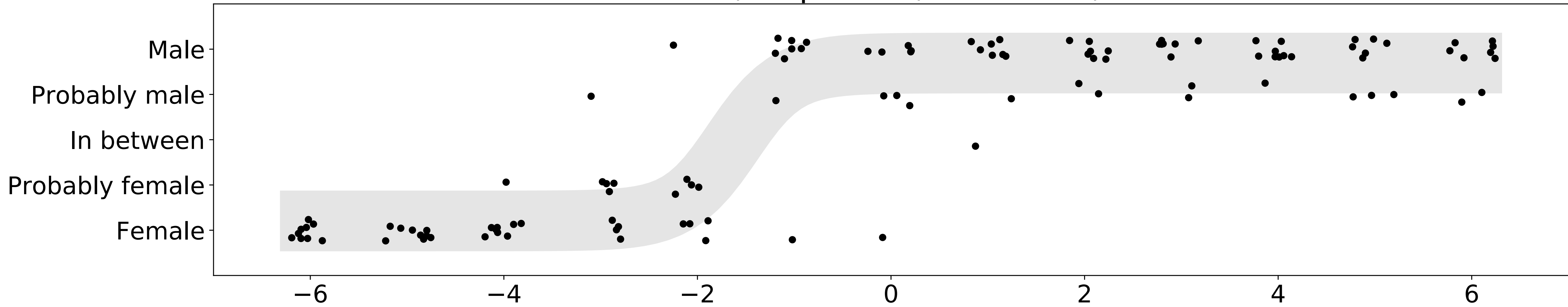
fit: shift=-3.0, slope=0.29, min=0.82, max=2.3



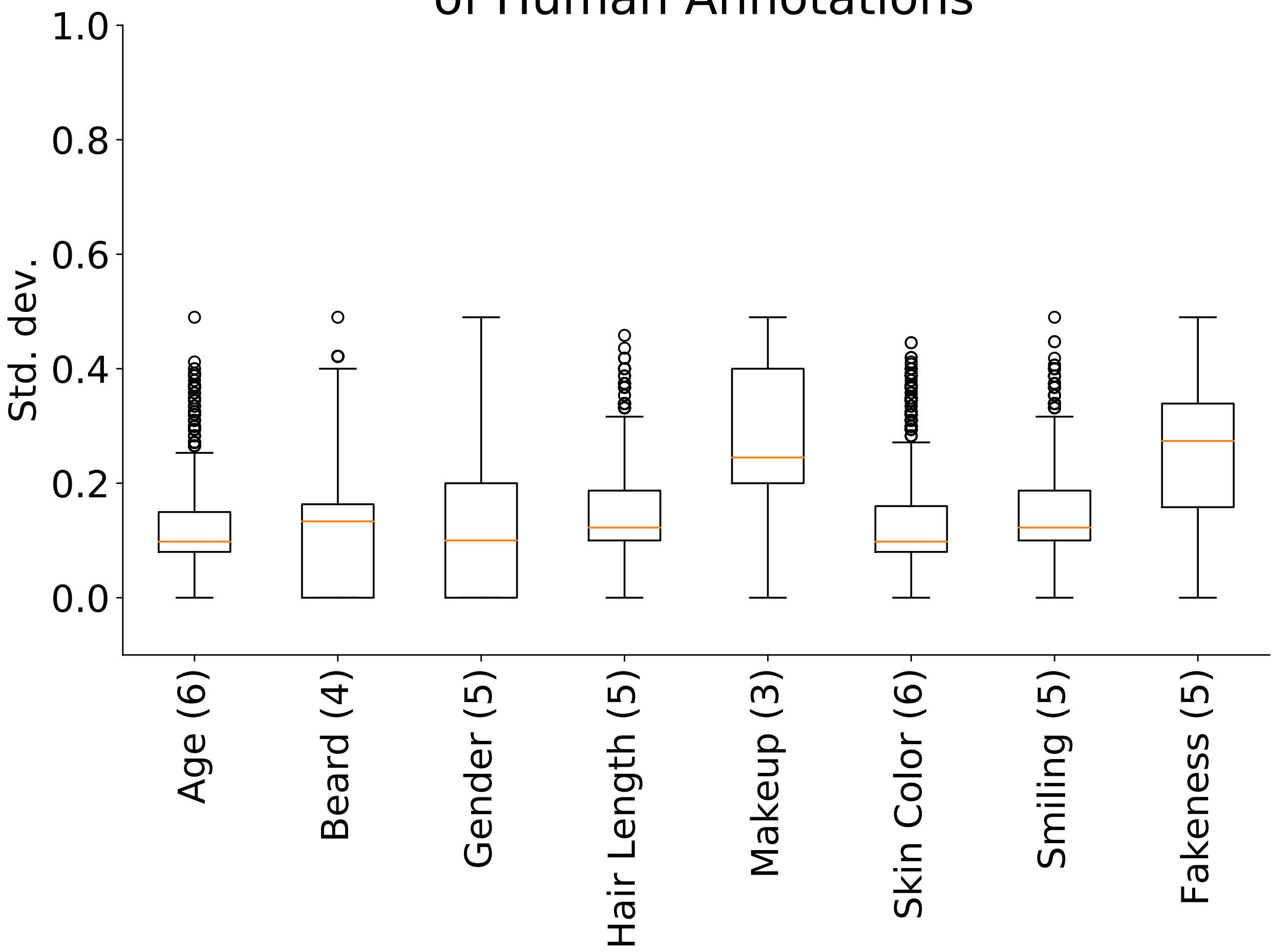
Gender



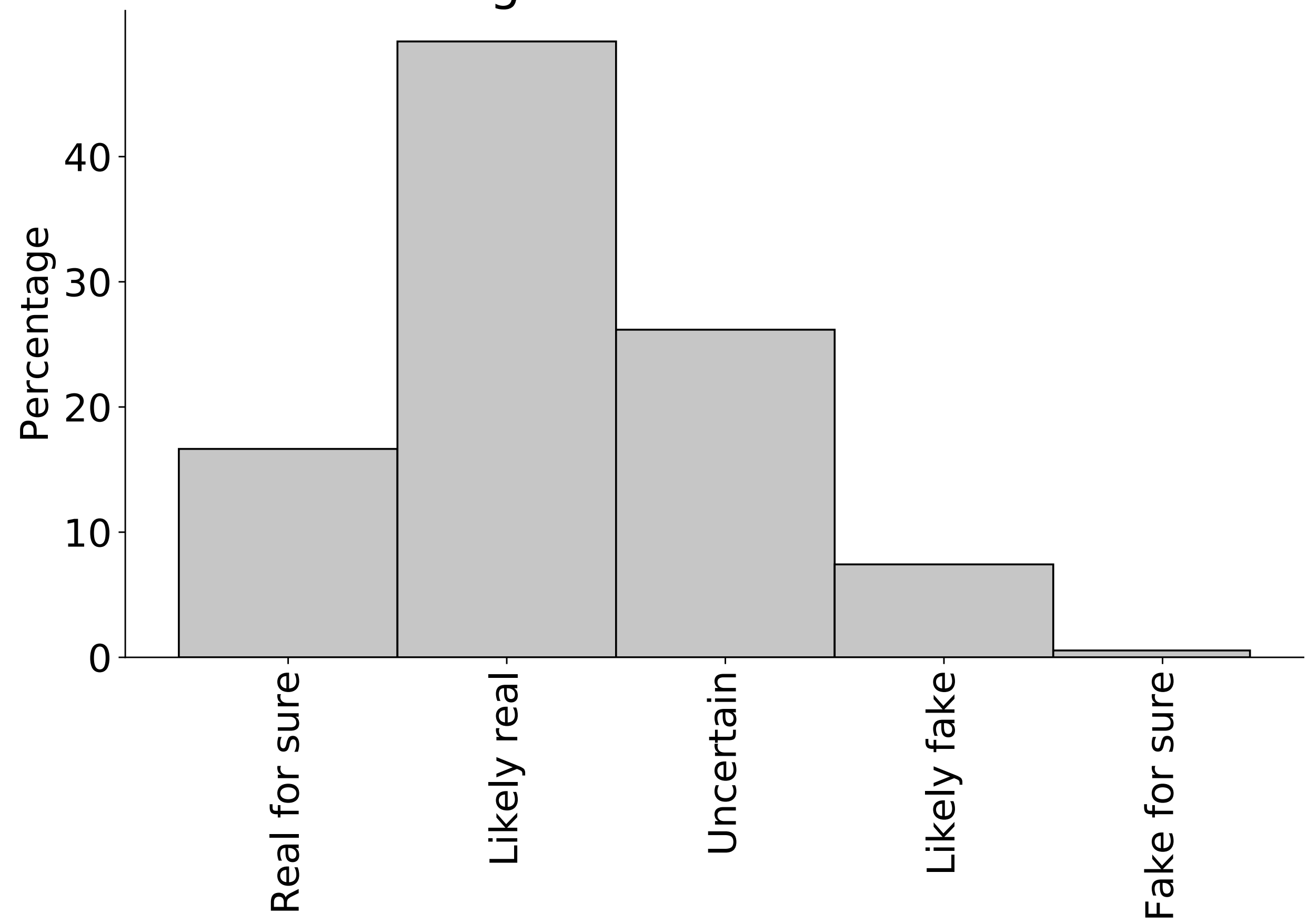
fit: shift=-1.6, slope=0.9, min=0.21, max=3.7



Per-Image Standard Deviations of Human Annotations



Per-Image Mean Fakeness Scores



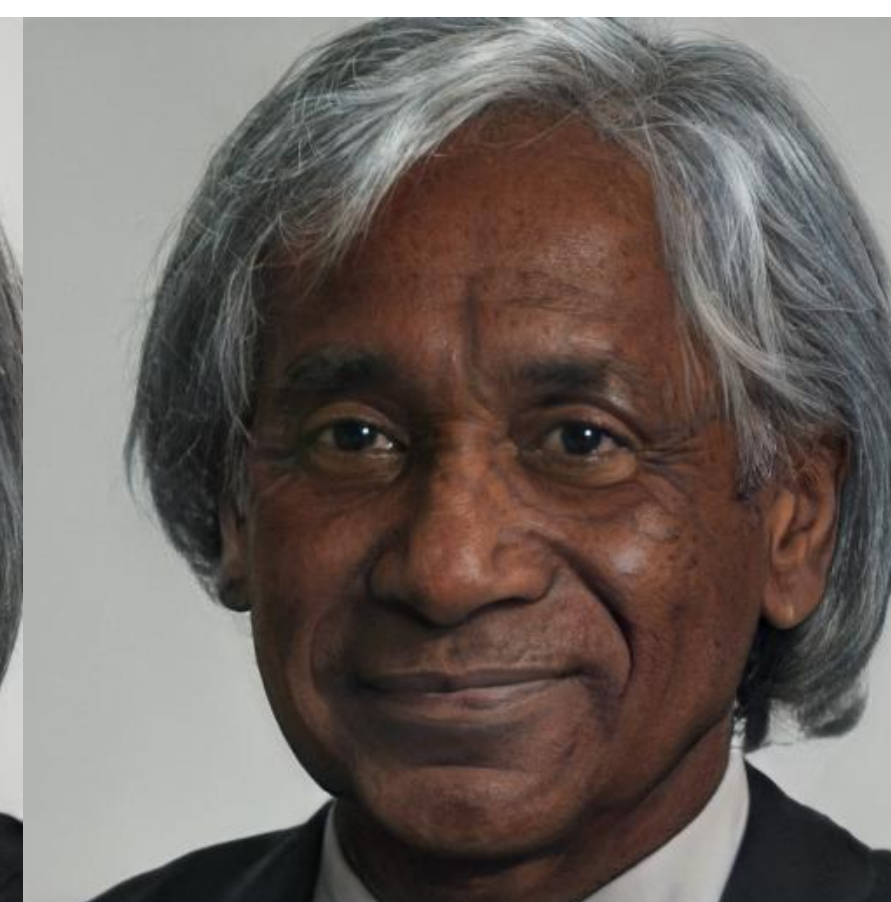


Humans: **Male**
 Intended: **Female**

Male
 Male

Likely Female
 Female

Male
 Male



Likely Female
 Female

Male
 Male ⁷⁰

Female
 Female

Likely Male
 Male

Hair length

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



Gender

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



Skin tone

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



Gender

< 0.2

0.2-0.3

0.3-0.4

0.4-0.5

0.5-0.6

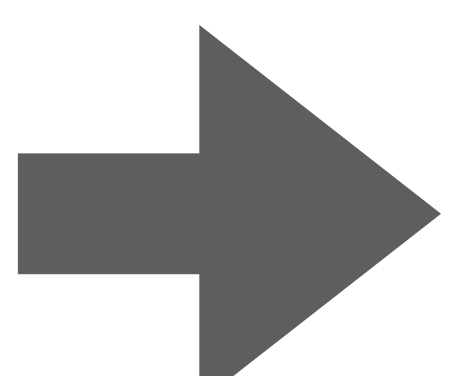
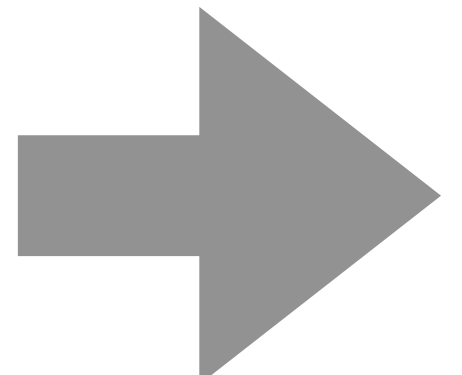
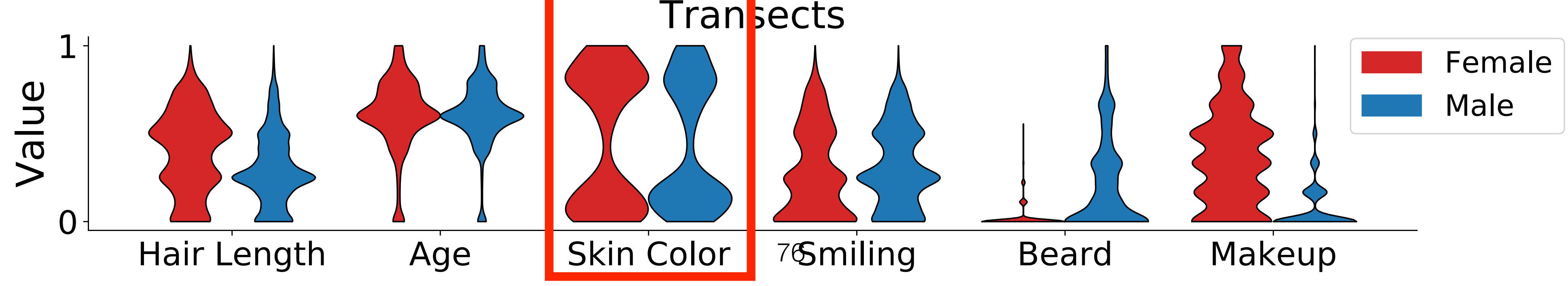
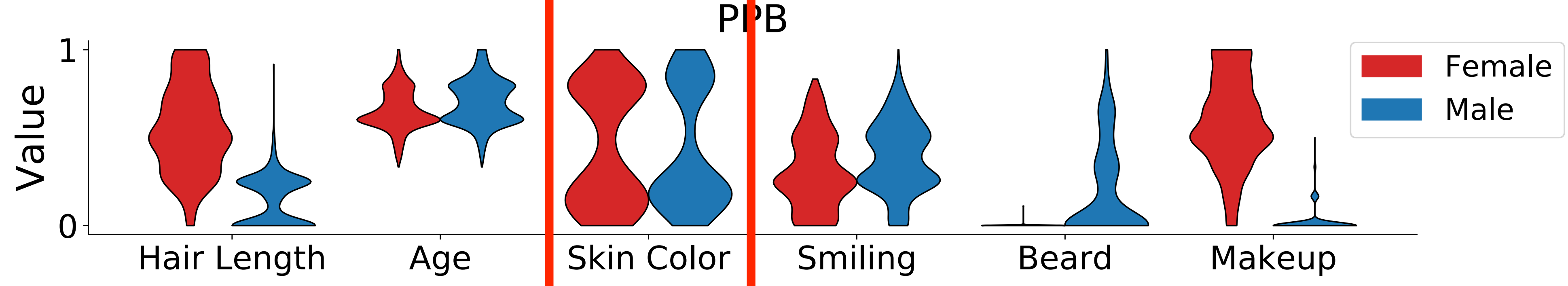
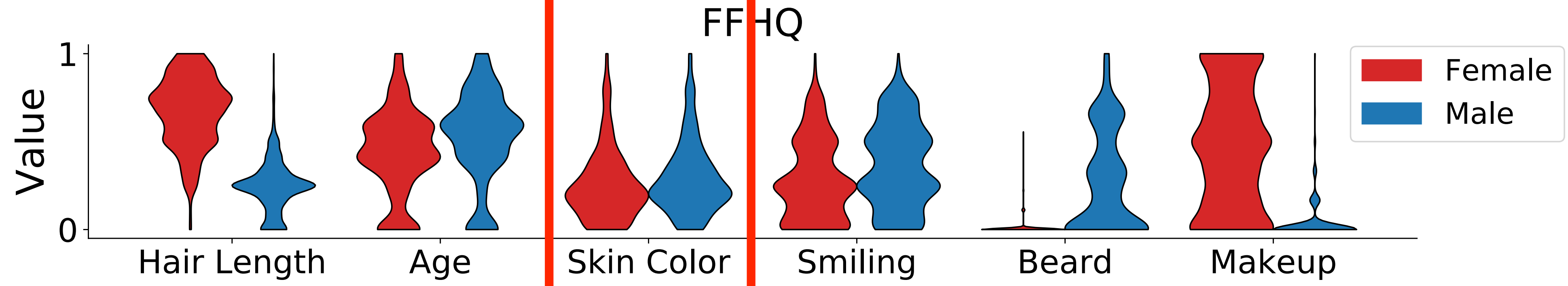
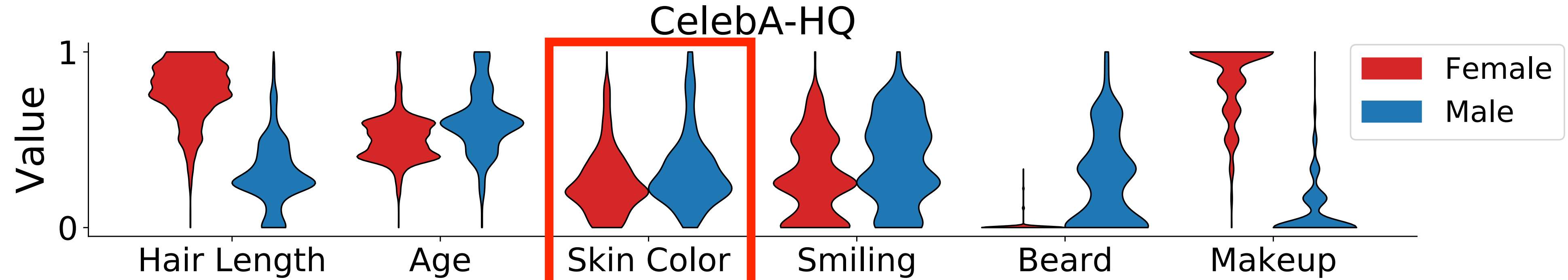
0.6-0.7

0.7-0.8

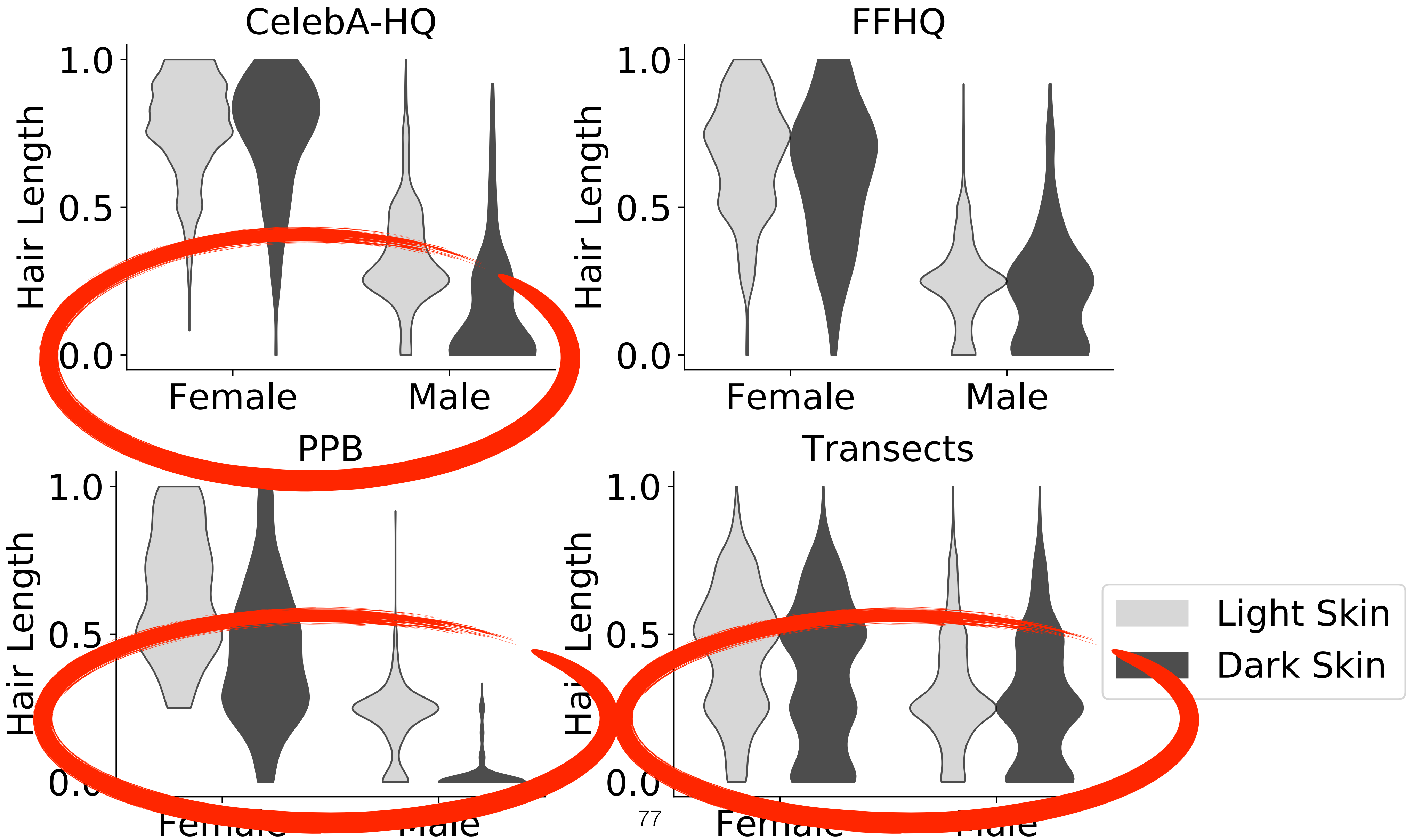
> 0.8



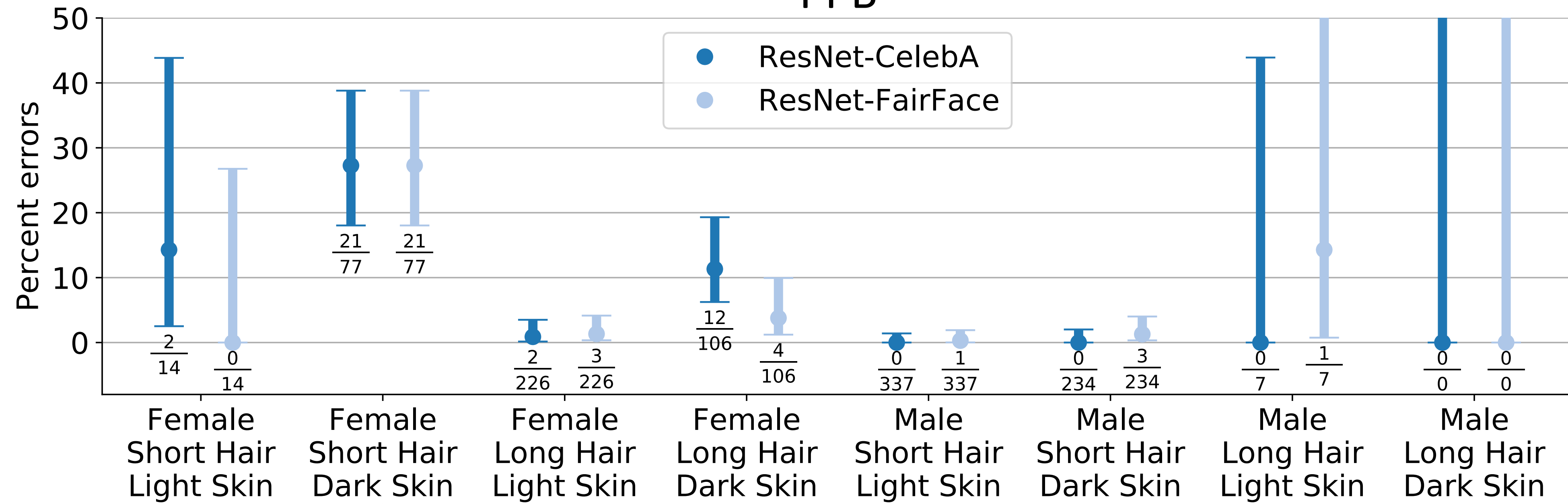
Data statistics and analysis



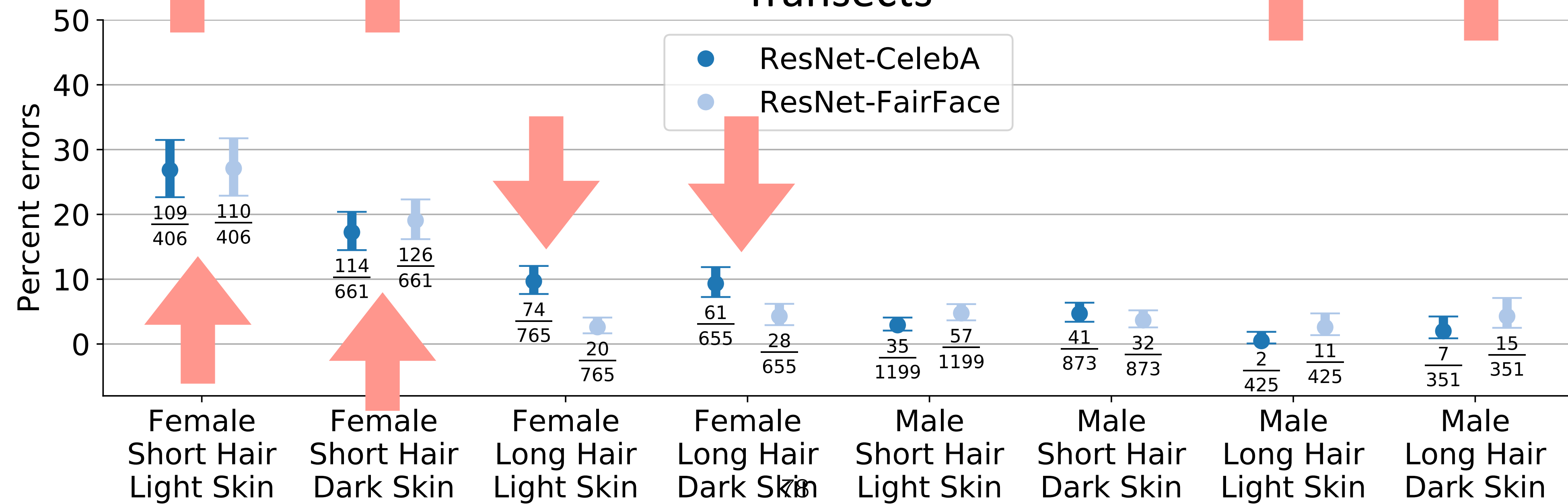
Intersectional

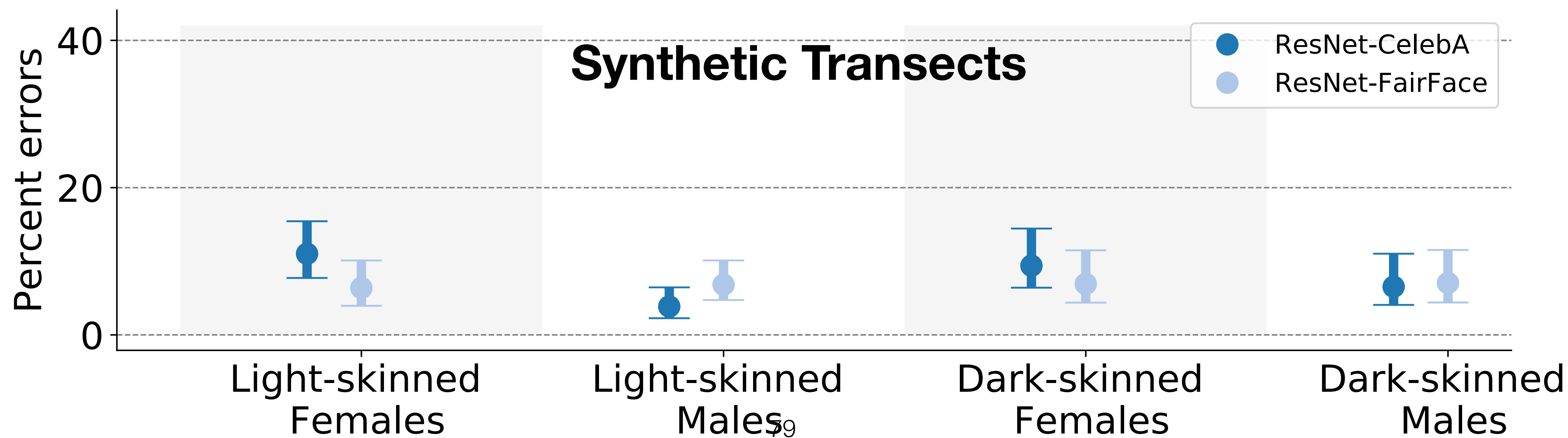
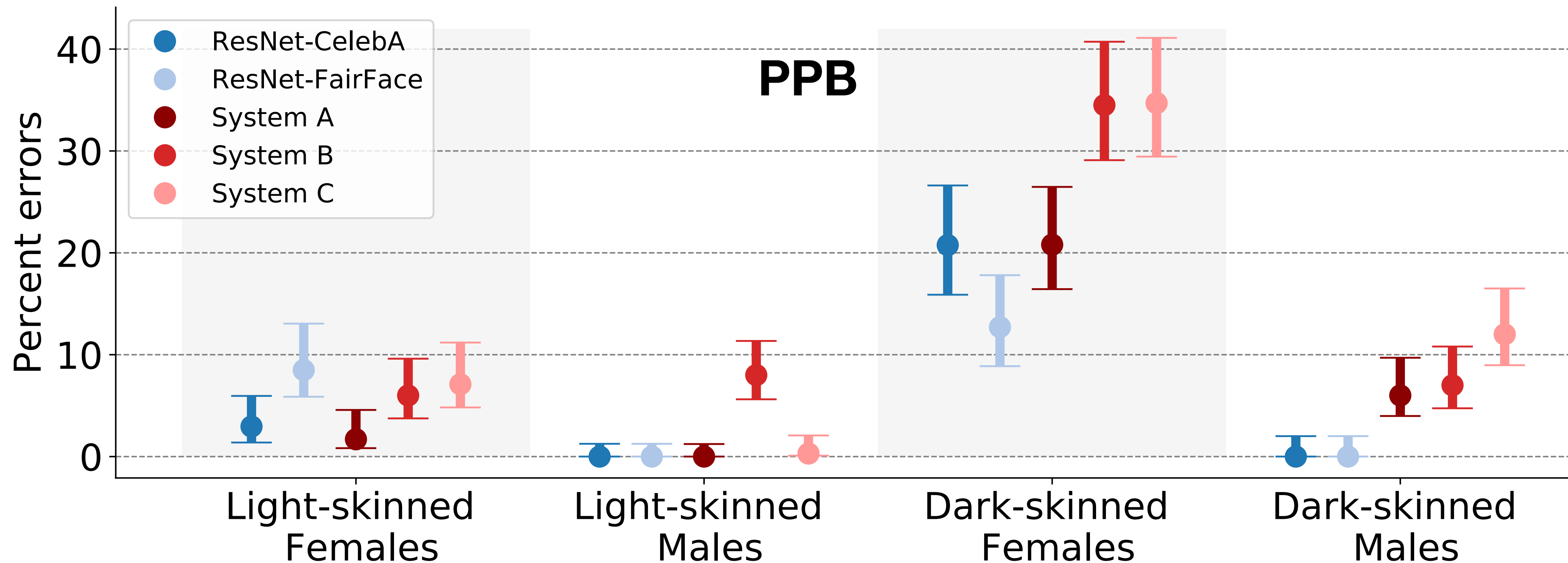


PPB



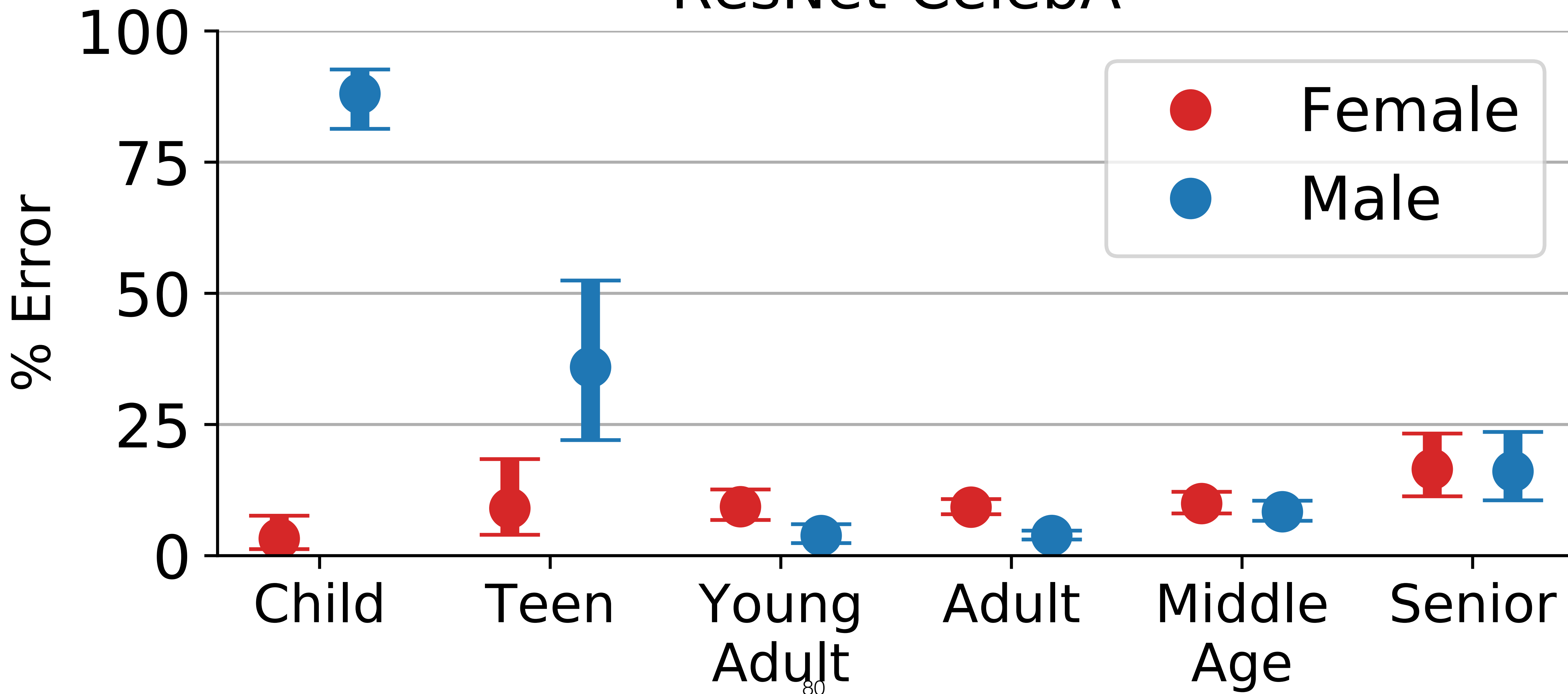
Transects





Age and gender

ResNet-CelebA



Residual correlations

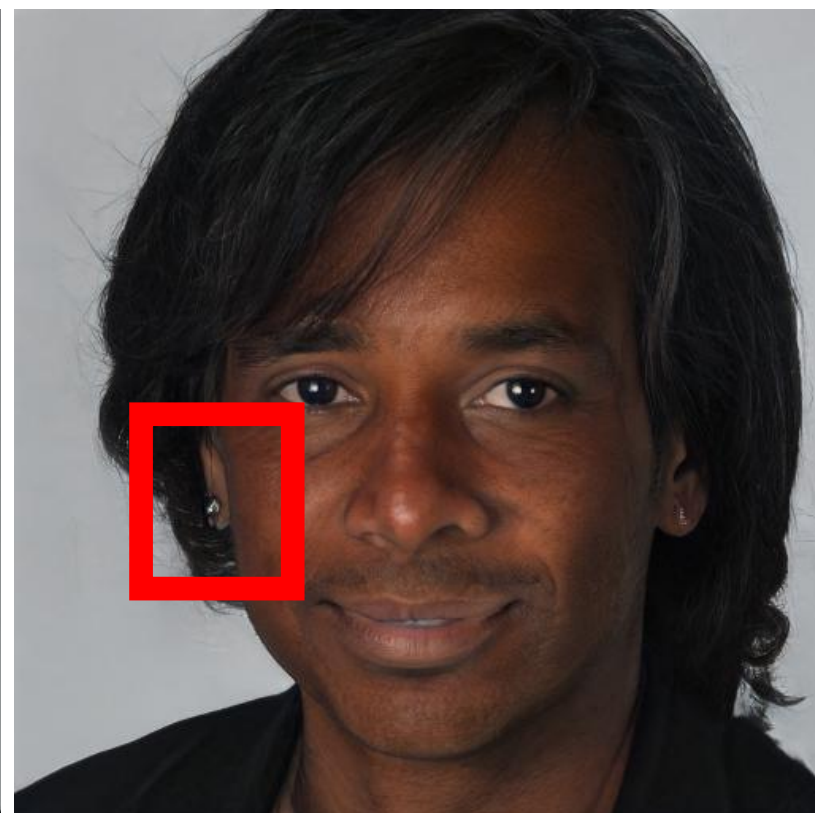
Example 1

Example 2

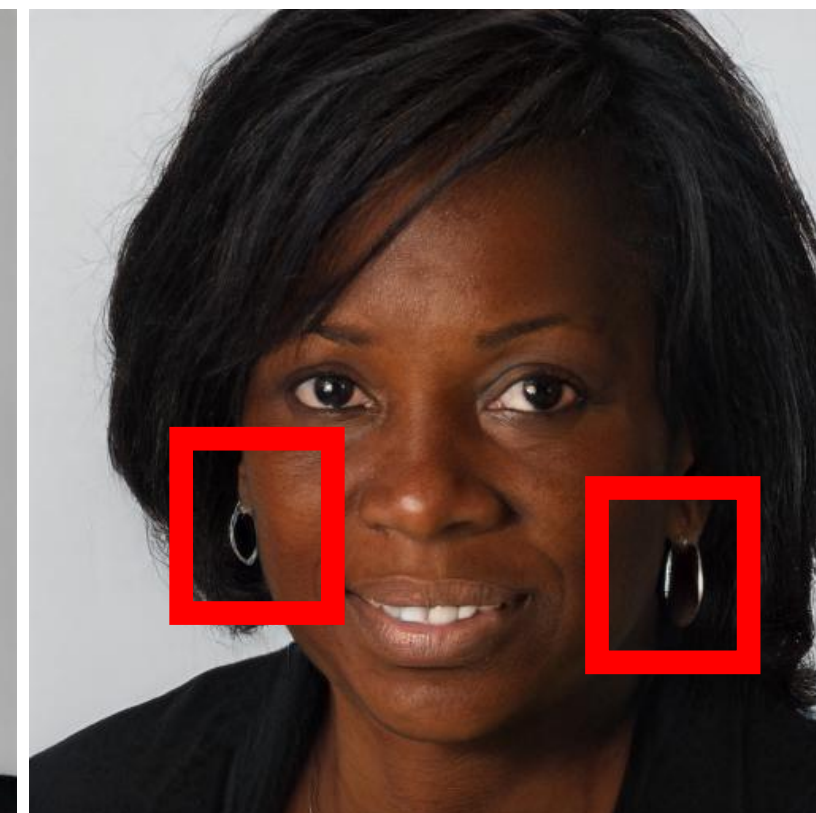
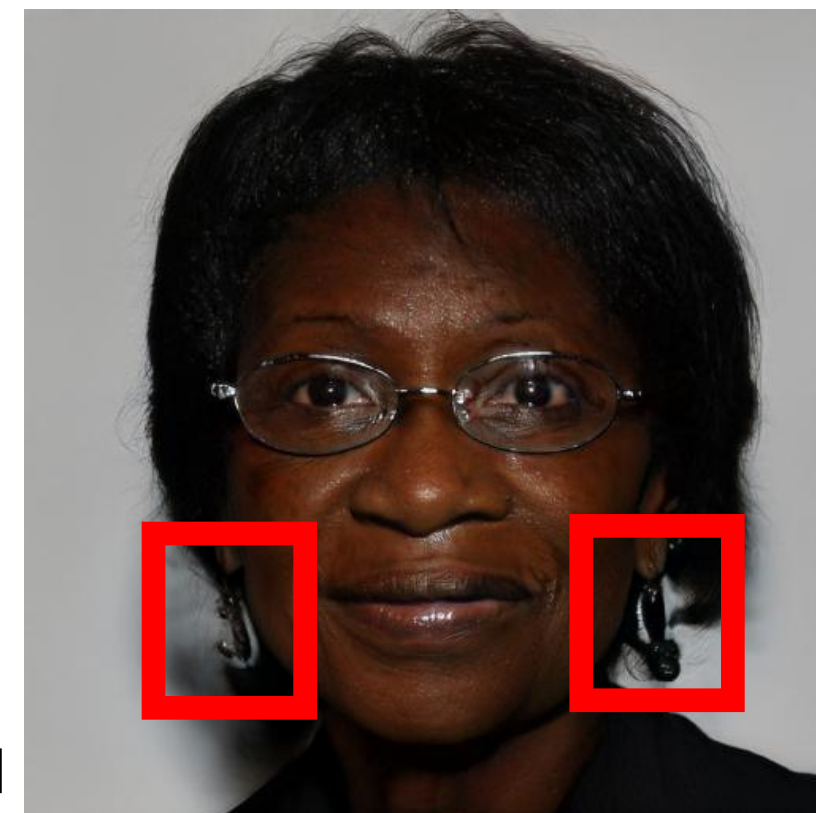


Males

Females



81



Quiz

Race Bias

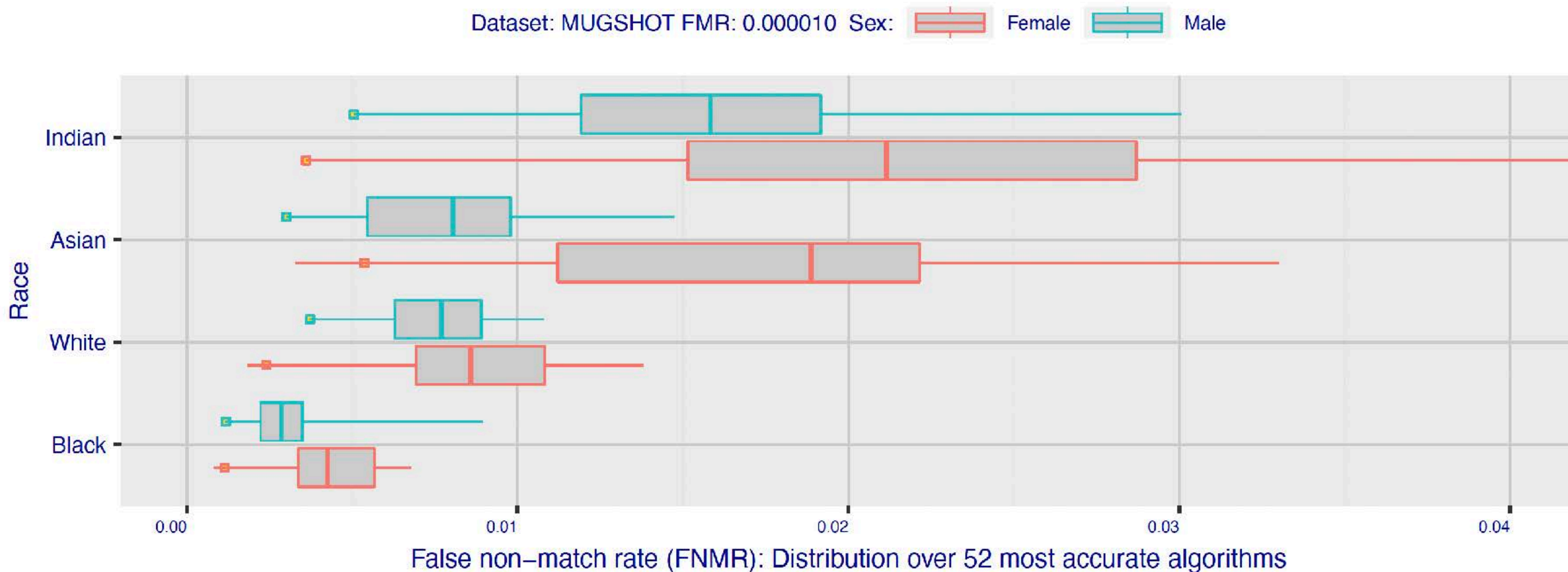


Figure 17: For mugshot comparisons, the figure shows the distribution of FNMR values over the 52 most accurate verification algorithms, by sex and race. The threshold was set for each algorithm to achieve $FMR = 0.00001$ over all imposter comparisons. The line within each box is the median over those algorithms; the box itself spans the interquartile range (26 algorithms) and the lines here extend to minimum and maximum values. The small box on the left side indicates the accuracy for best algorithm overall, on this dataset `alphaface801`.

Race bias

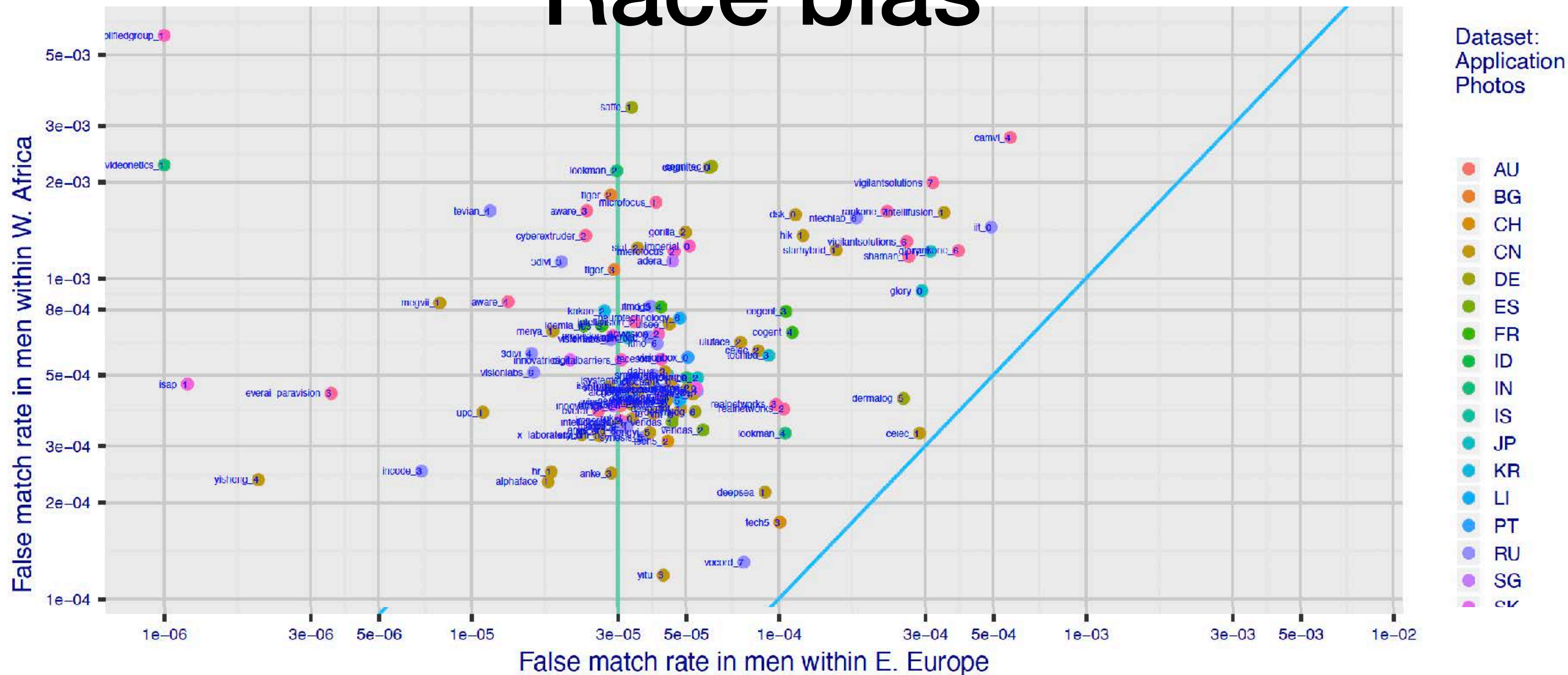
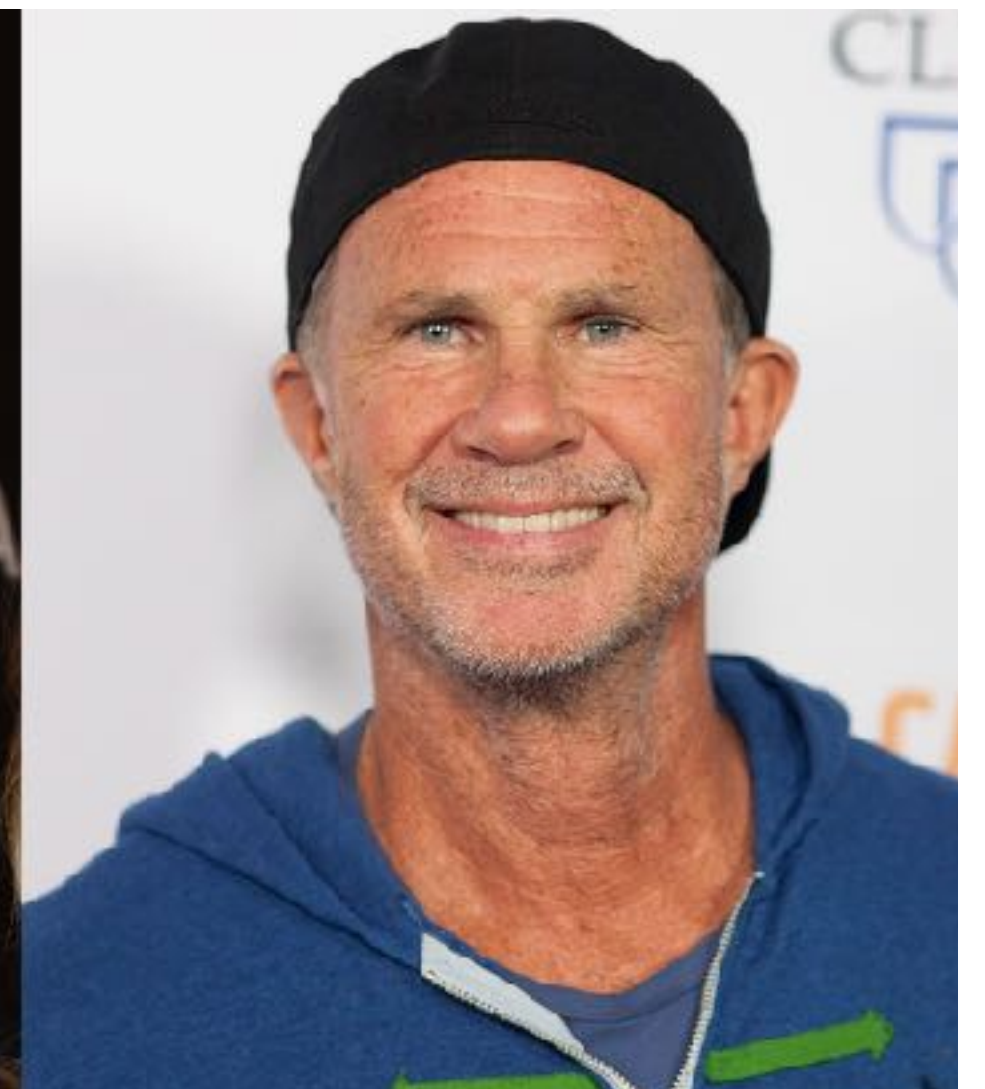
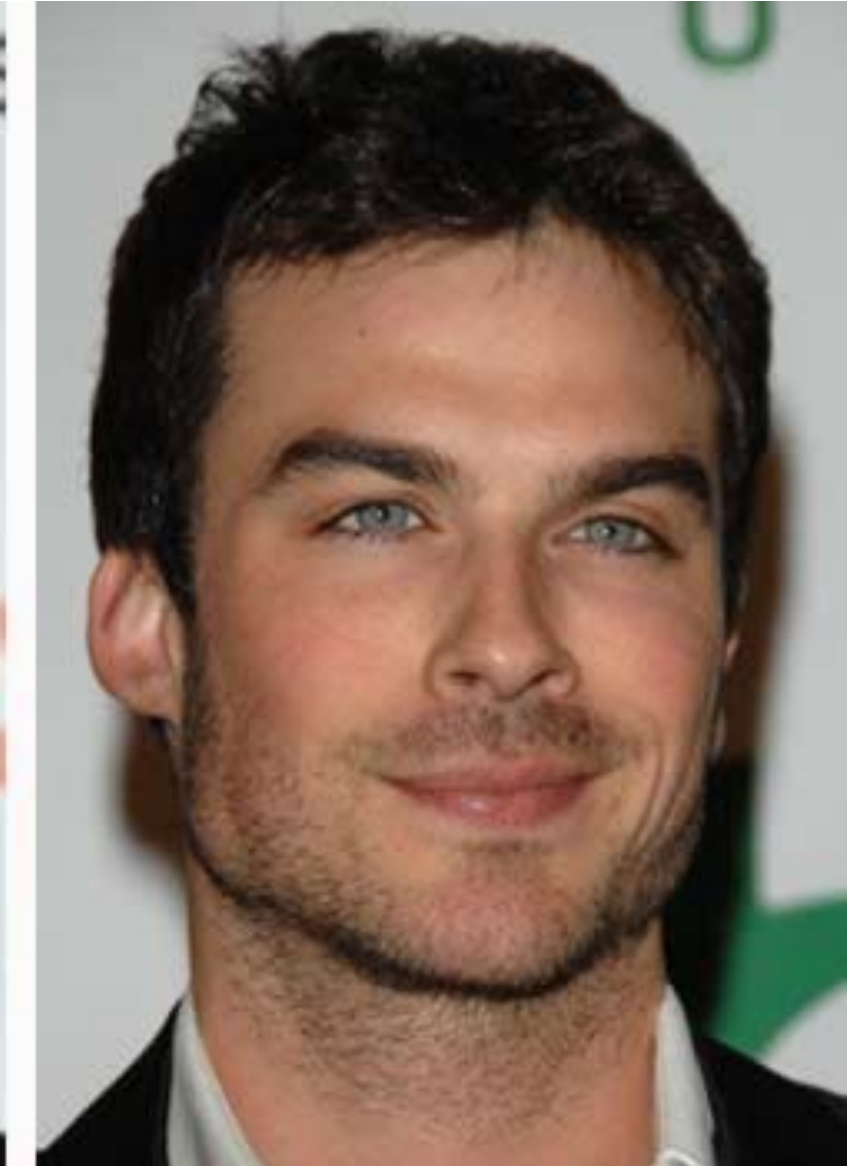


Figure 8: The scatter plot shows FMR when comparing same-age men within and across three Eastern European countries (Russia, Ukraine, Poland), against FMR obtained comparing men within and across three West African countries (Ghana, Liberia, Nigeria). The threshold is fixed for each algorithm to give the FMR noted in the annotation over white men in the U.S. mugshot database. This is indicated by the vertical and horizontal green lines. The blue diagonal line $y = x$ is included to show “one-to-one” comparisons. The color coding identifies the country of the algorithm developer.

Identity



WM

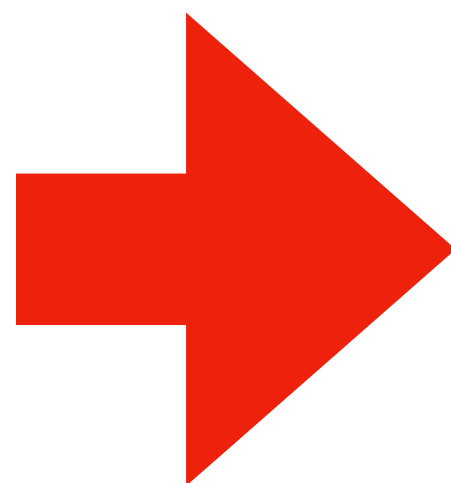
WF

BM

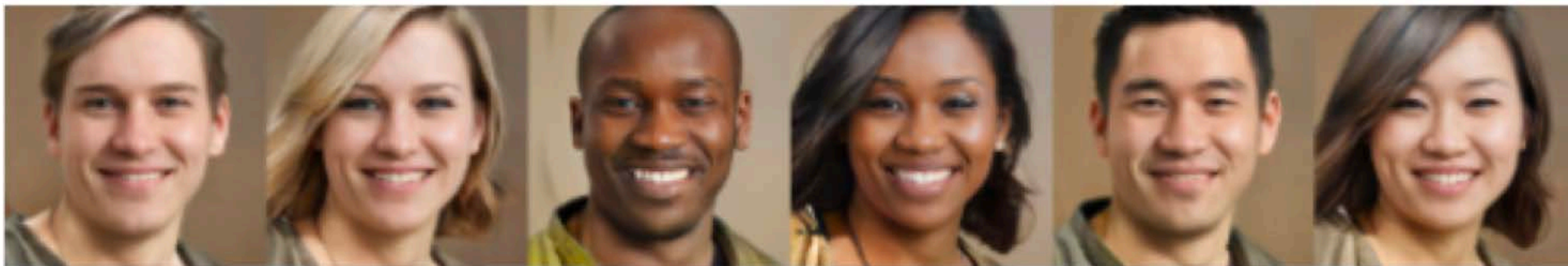
BF

AM

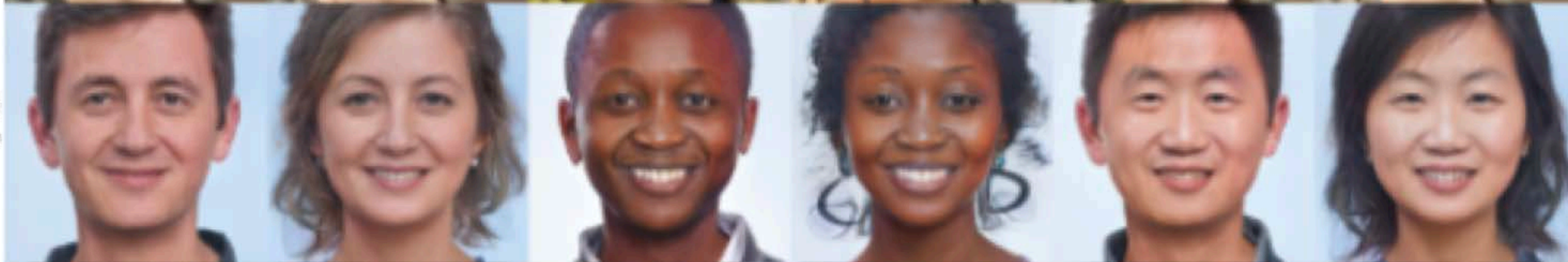
AF



S1



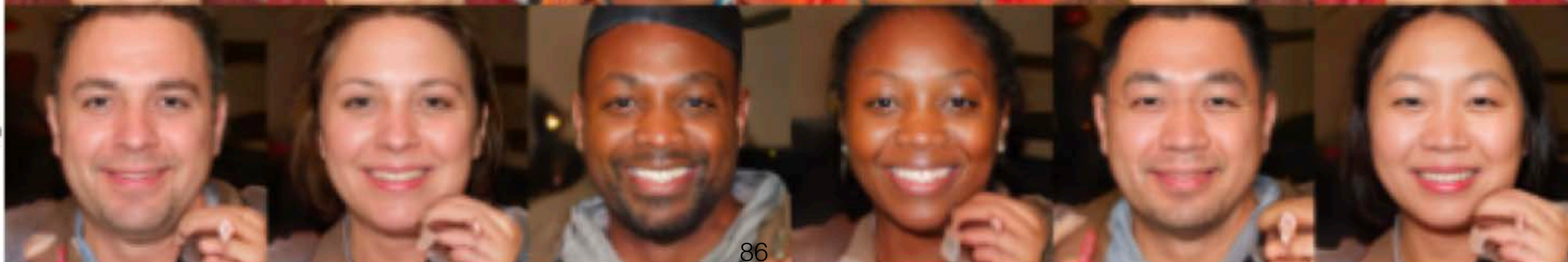
S2



S3



S4



Pose



Lighting

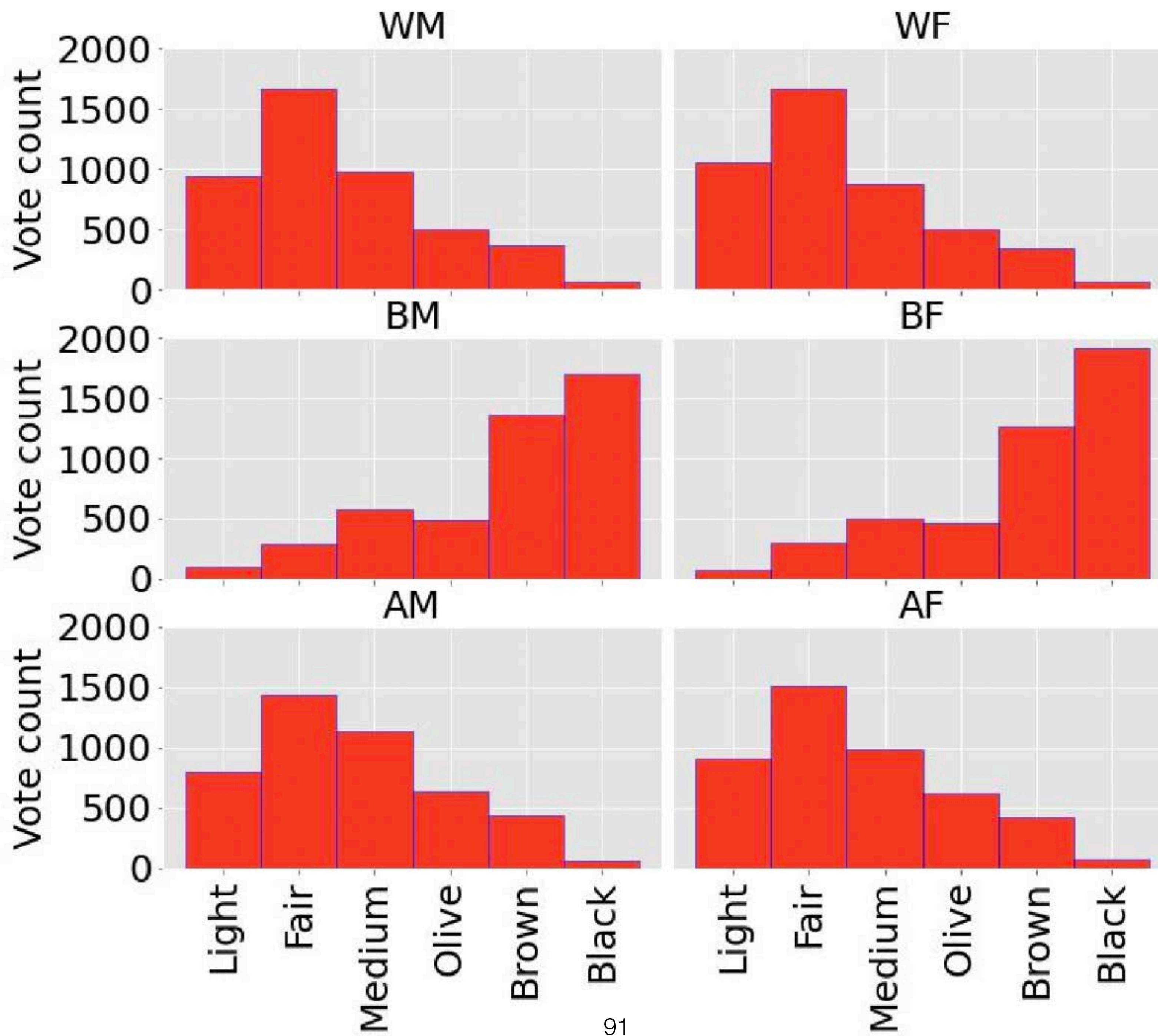


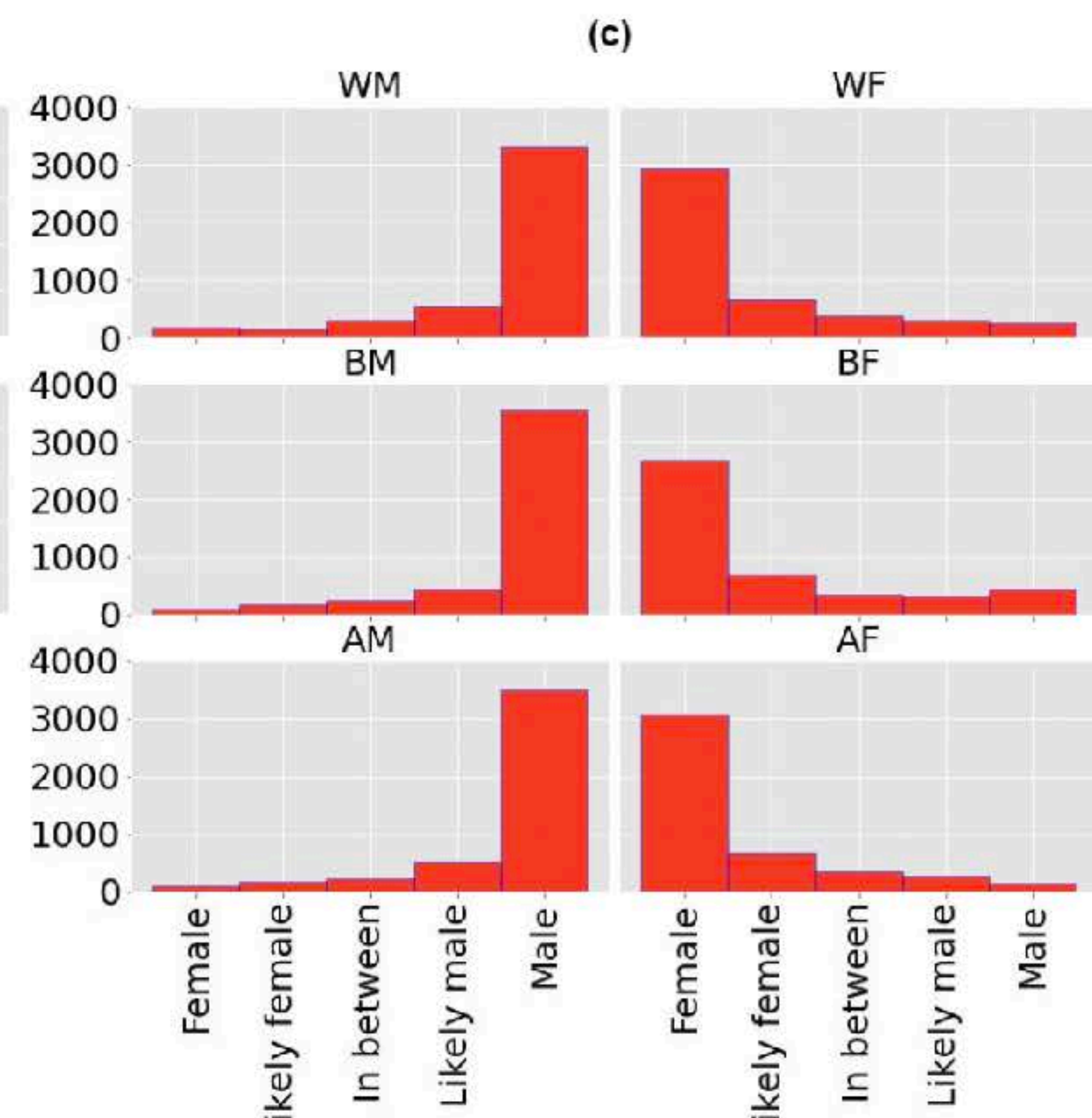
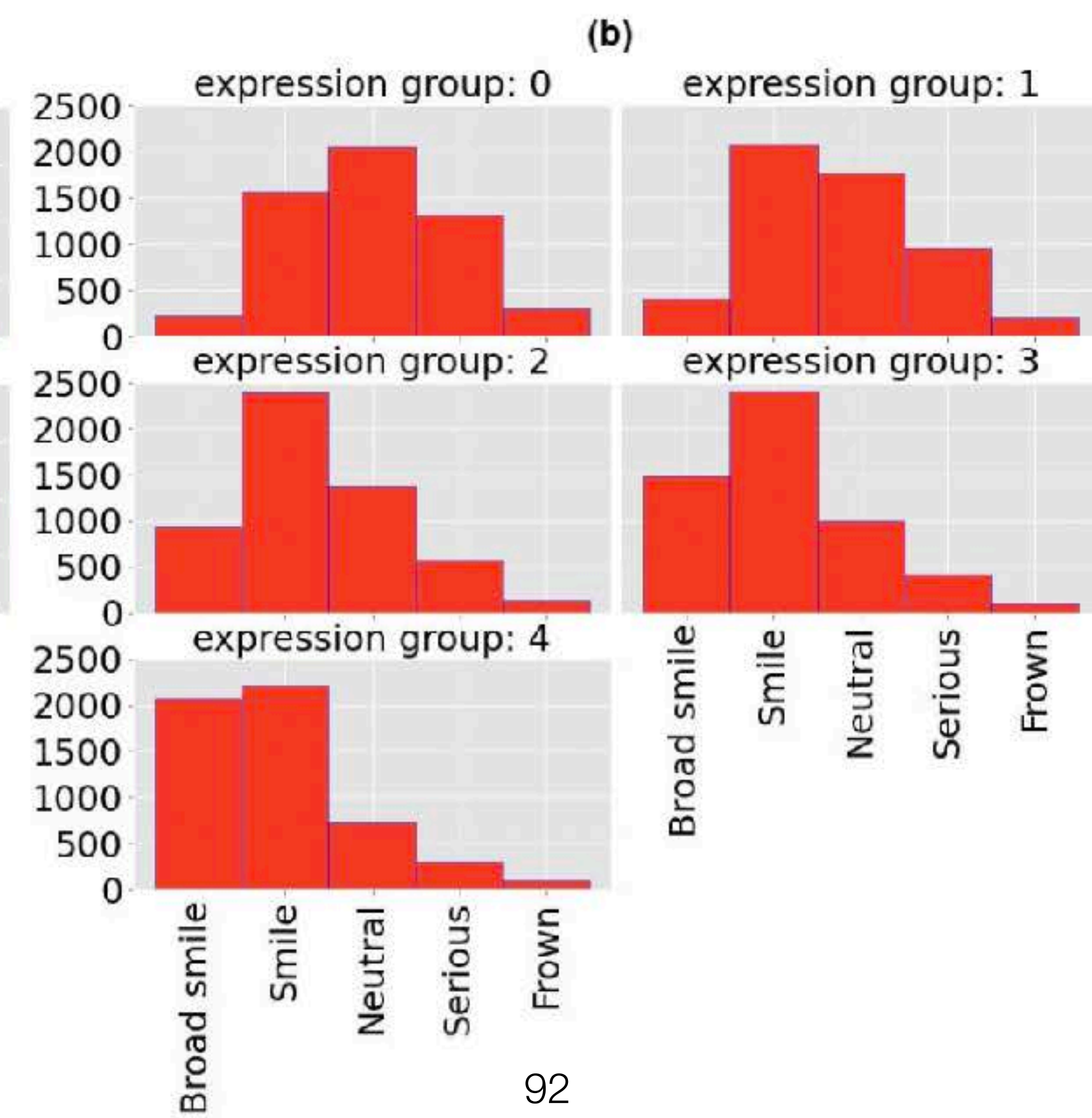
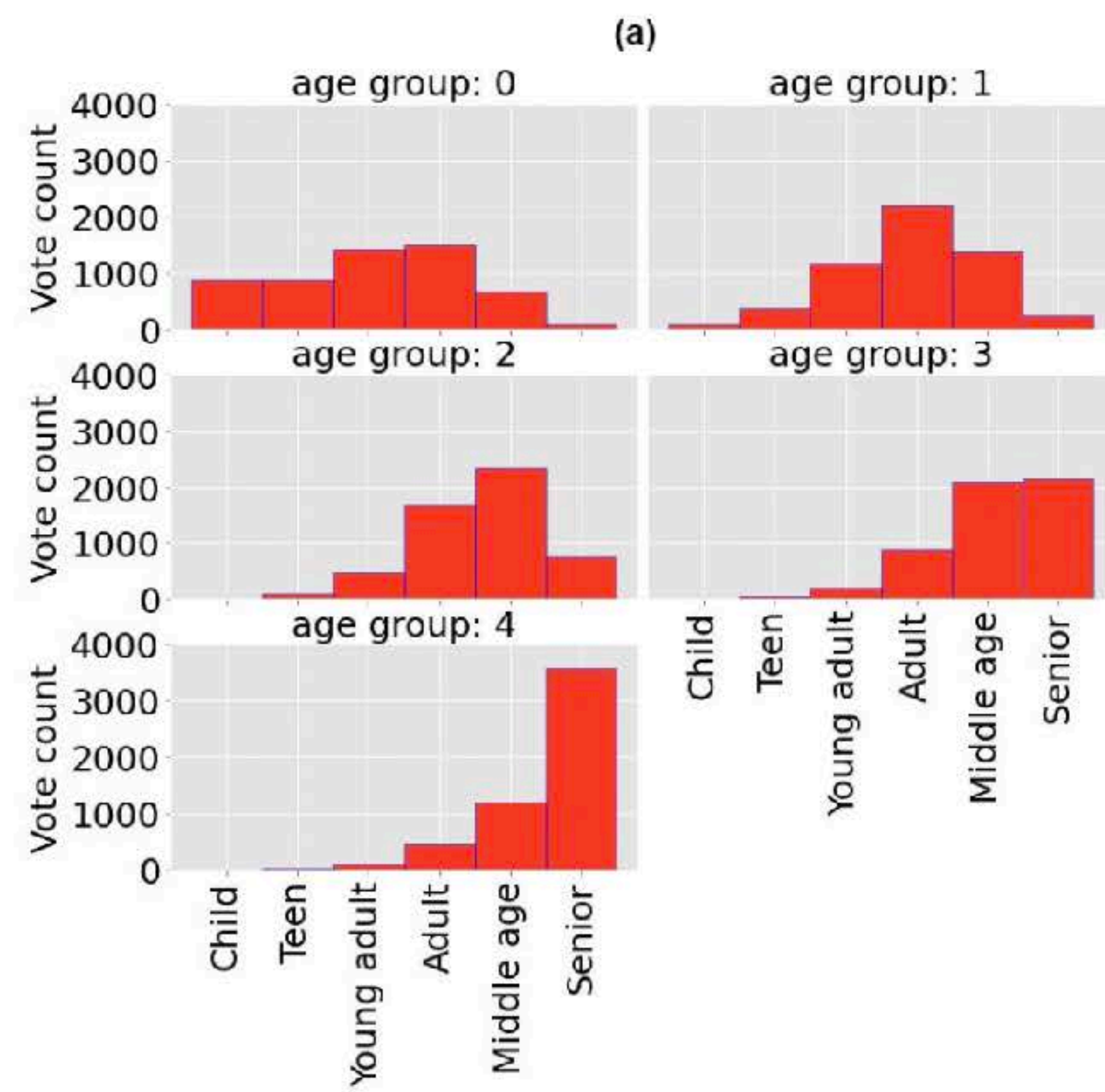
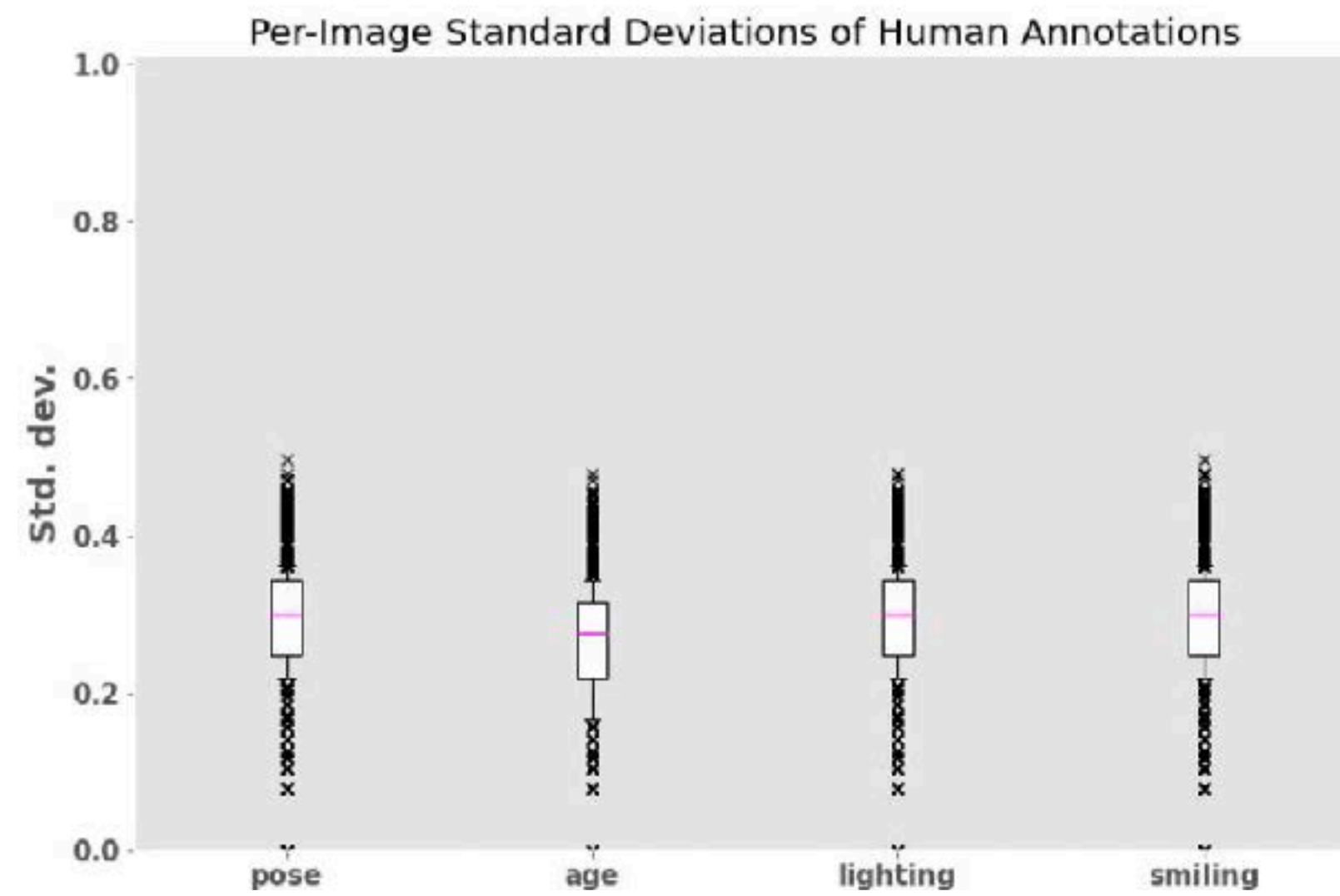
Age



Expression







WM

WF

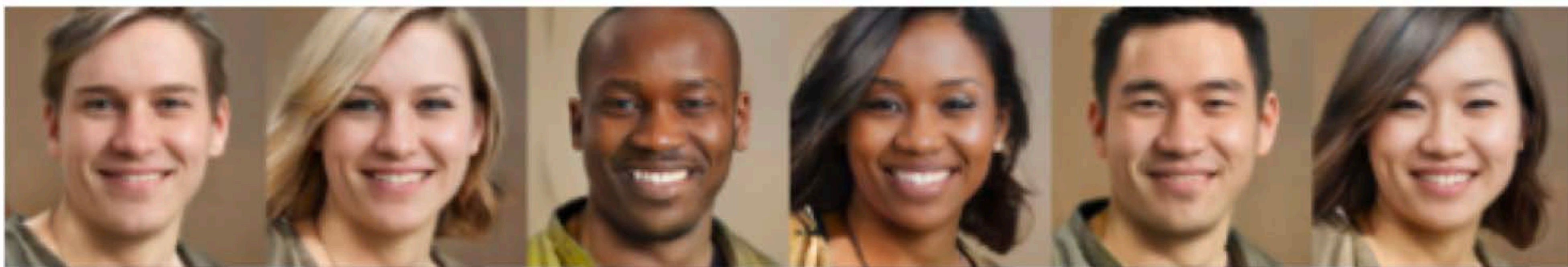
BM

BF

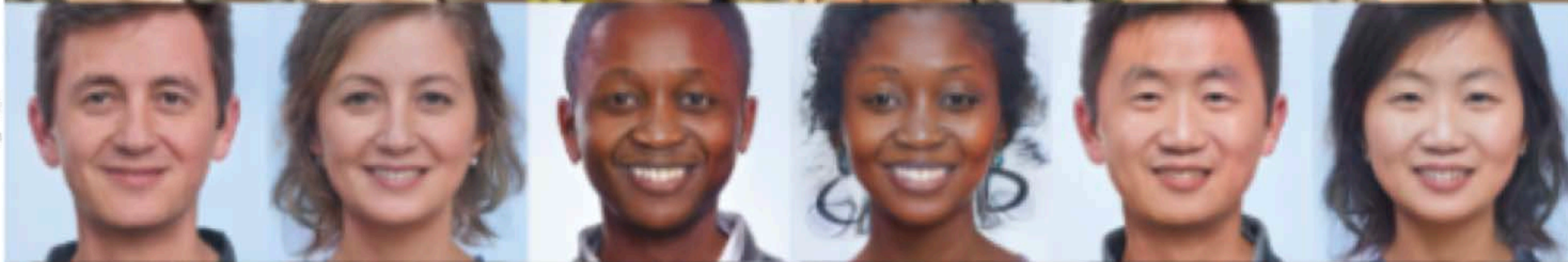
AM

AF

S1



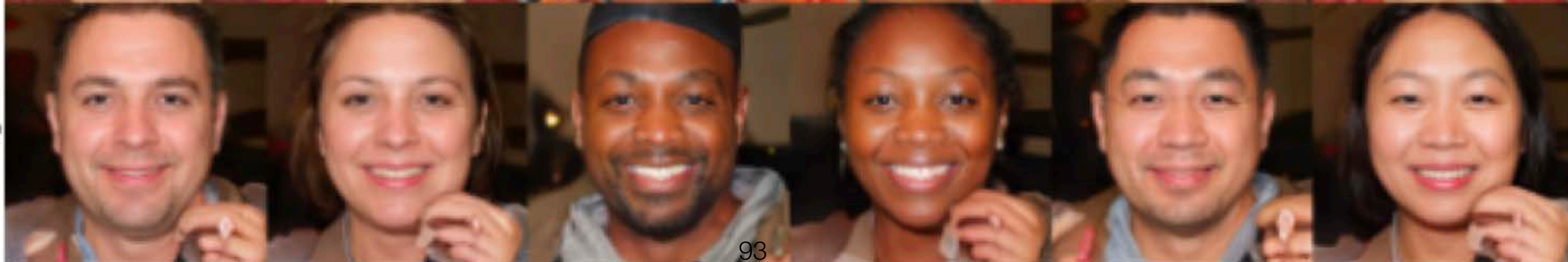
S2

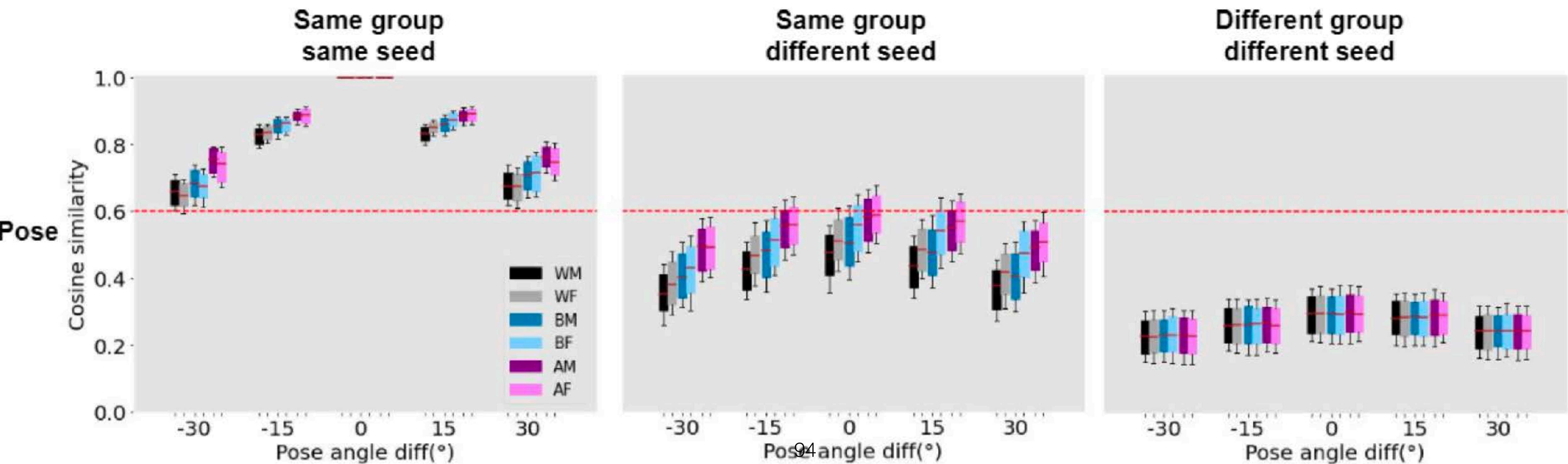


S3



S4





Instructions

Shortcuts

Same or different person? Are the two people the same (OK different age, hair style, facial expression), or different (maybe a close relative dressed and styled differently)?

Instructions



Likely same - Most people would think that they are the same.

Possibly same - Many but not all people would think that they are the same.

Not sure - People would be split on this.

Possibly different - Perhaps siblings. Many but not all people would think different.

Likely different - Most people would say different.



Select an option

Likely same 1

Possibly same 2

Not sure 3

Possibly different 4

Likely different 5

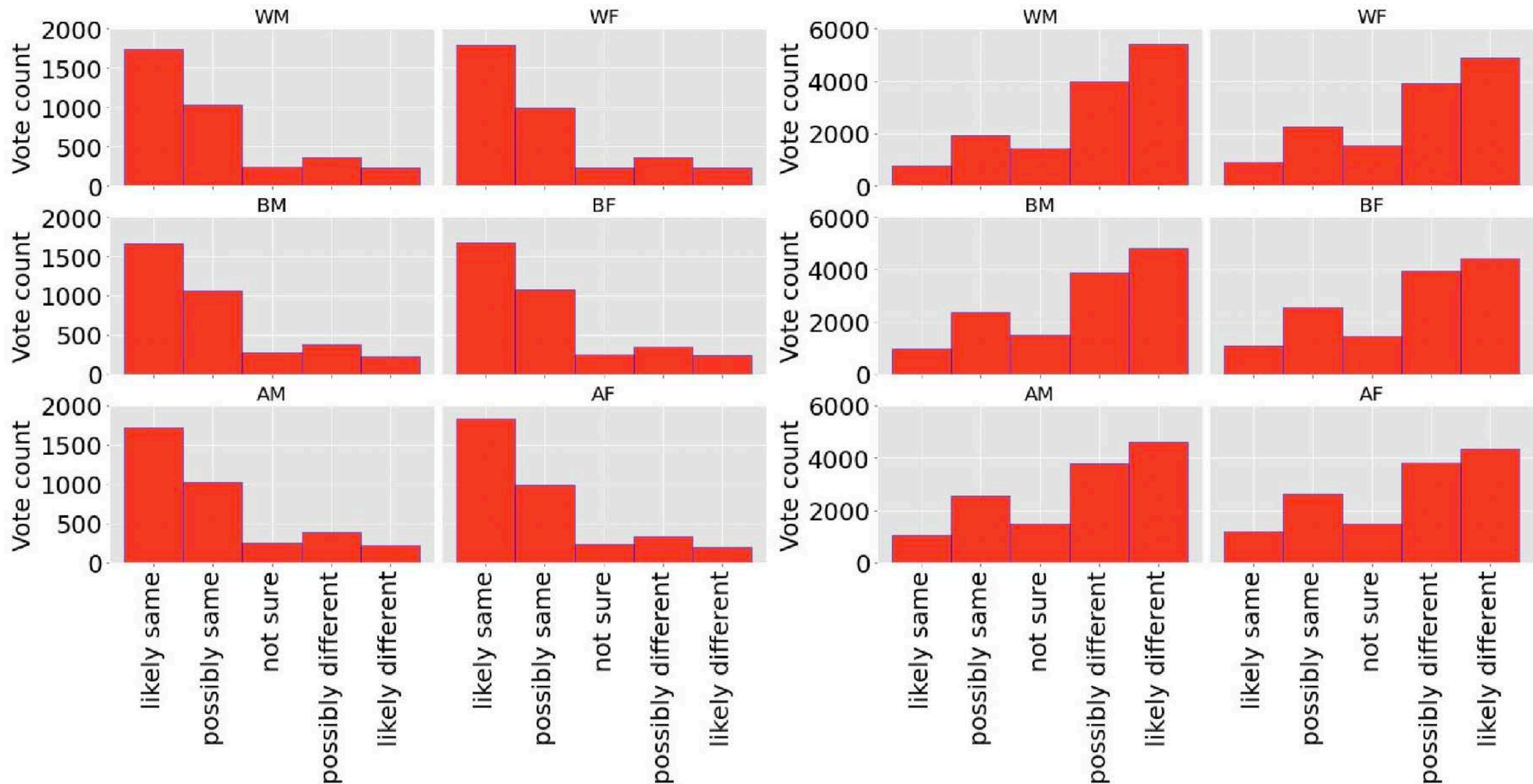
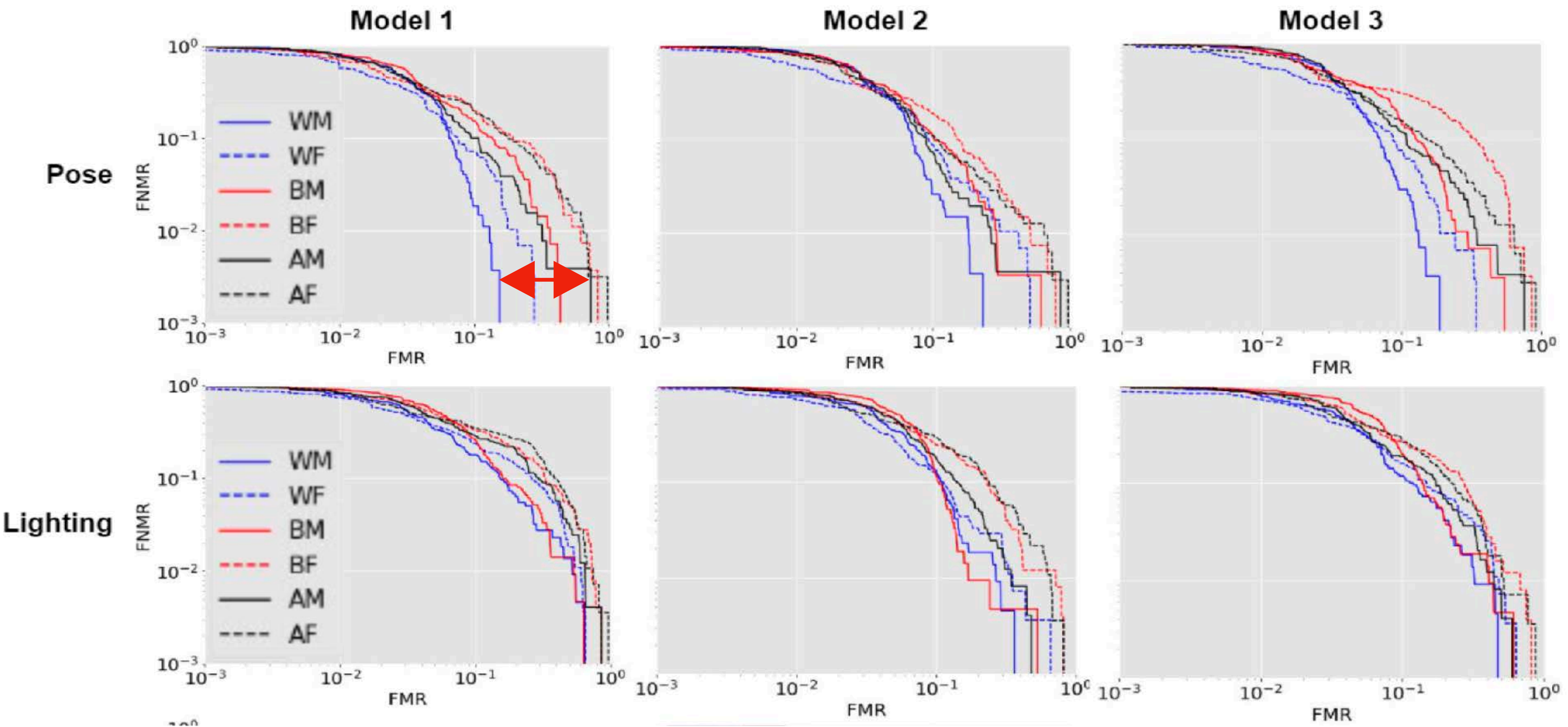




Figure 4. **Human annotation results on different pose image pairs.** (a): face images of different poses from the same seed and same demographic group. (b): face images of different poses from different seed but same demographic group. This trend is consistent across attributes, see Supplementary for results on other attributes.



Towards *responsible* responsible AI

- Observational methods 
- Experimental method 
- Synthetic samples + human annotations

Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale

Federico Bianchi^{*1}, Pratyusha Kalluri^{*1}, Esin Durmus^{*1}, Faisal Ladhak^{*1,2}, Myra Cheng^{*1},
Debora Nozza³, Tatsunori Hashimoto¹, Dan Jurafsky^{†1}, James Zou^{†1}, and Aylin Caliskan^{†4}

¹Stanford University

²Columbia University

³Bocconi University

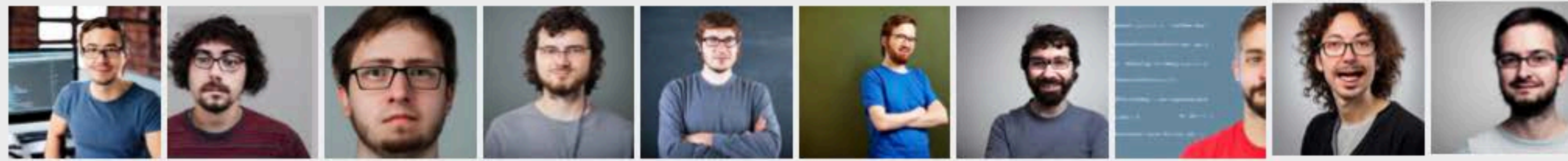
⁴University of Washington

Abstract

Machine learning models are now able to convert user-written text descriptions into naturalistic images. These models are available to anyone online and are being used to generate millions of images a day. We investigate these models and find that they amplify dangerous and complex stereotypes. Moreover, we find that the amplified stereotypes are difficult to predict and not easily mitigated by users or model owners. The extent to which these image-generation models perpetuate and amplify stereotypes and their mass deployment is cause for serious concern.

1 Introduction

a software developer



—> Young, pale, male, facial hair, eyeglasses

a flight attendant



—> Female, young, pale

a chef



—> Male

a cook



—> Female, brown

a taxi driver



—> Male, brown

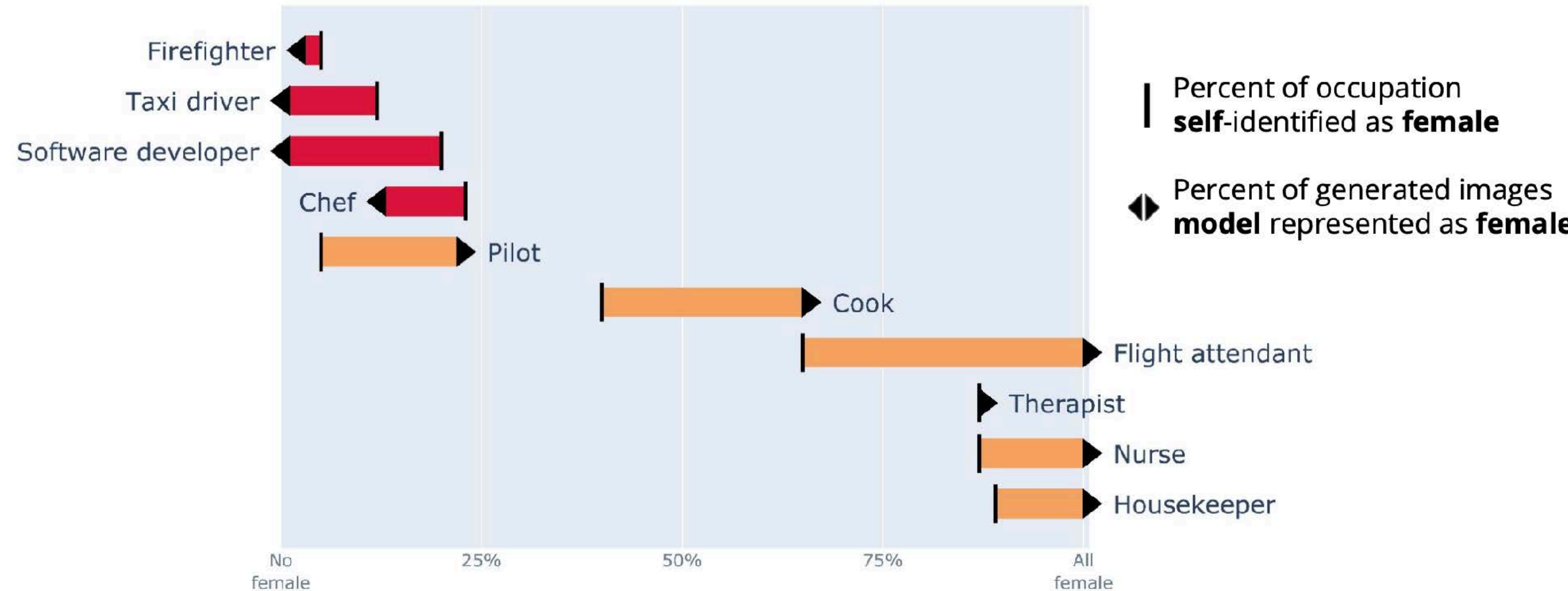
a housekeeper



—> Female, brown

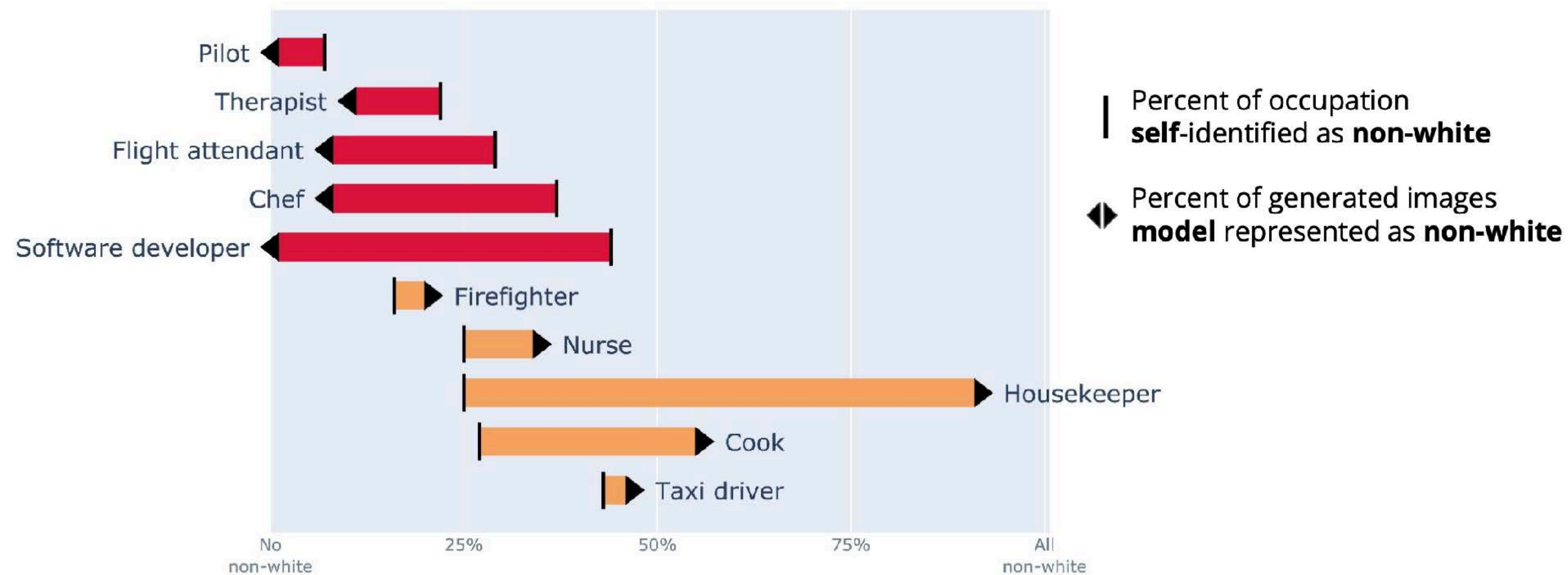
Comparison with bureau of labor data

Percent of occupation identified as female



Comparison with bureau of labor data

Percent of occupation identified as non-white



Improvability

The New York Times

ECONOMIC VIEW

Biased Algorithms Are Easier to Fix Than Biased People

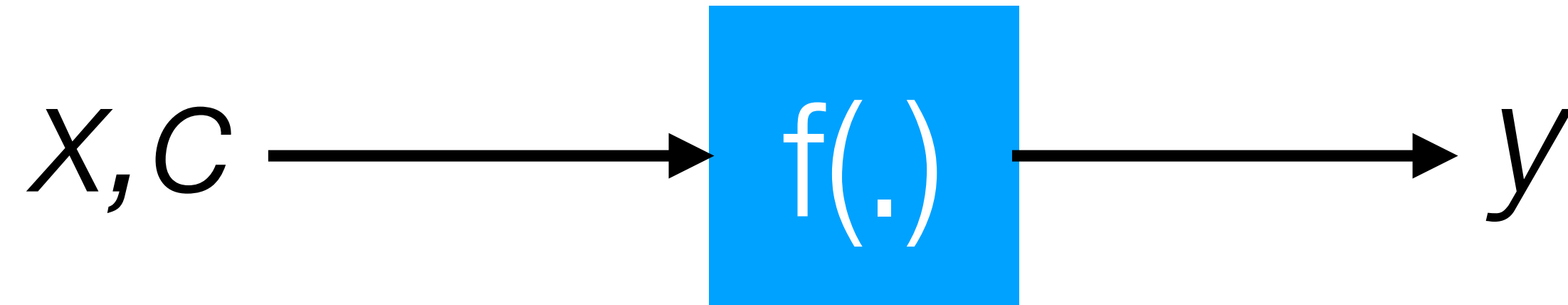
Racial discrimination by algorithms or by people is harmful — but that's where the similarities end.

By Sendhil Mullainathan

Dec. 6, 2019

Fixing bias?

Modify the loss function



$x = \text{action, state}$

$c = \text{protected attribute}$

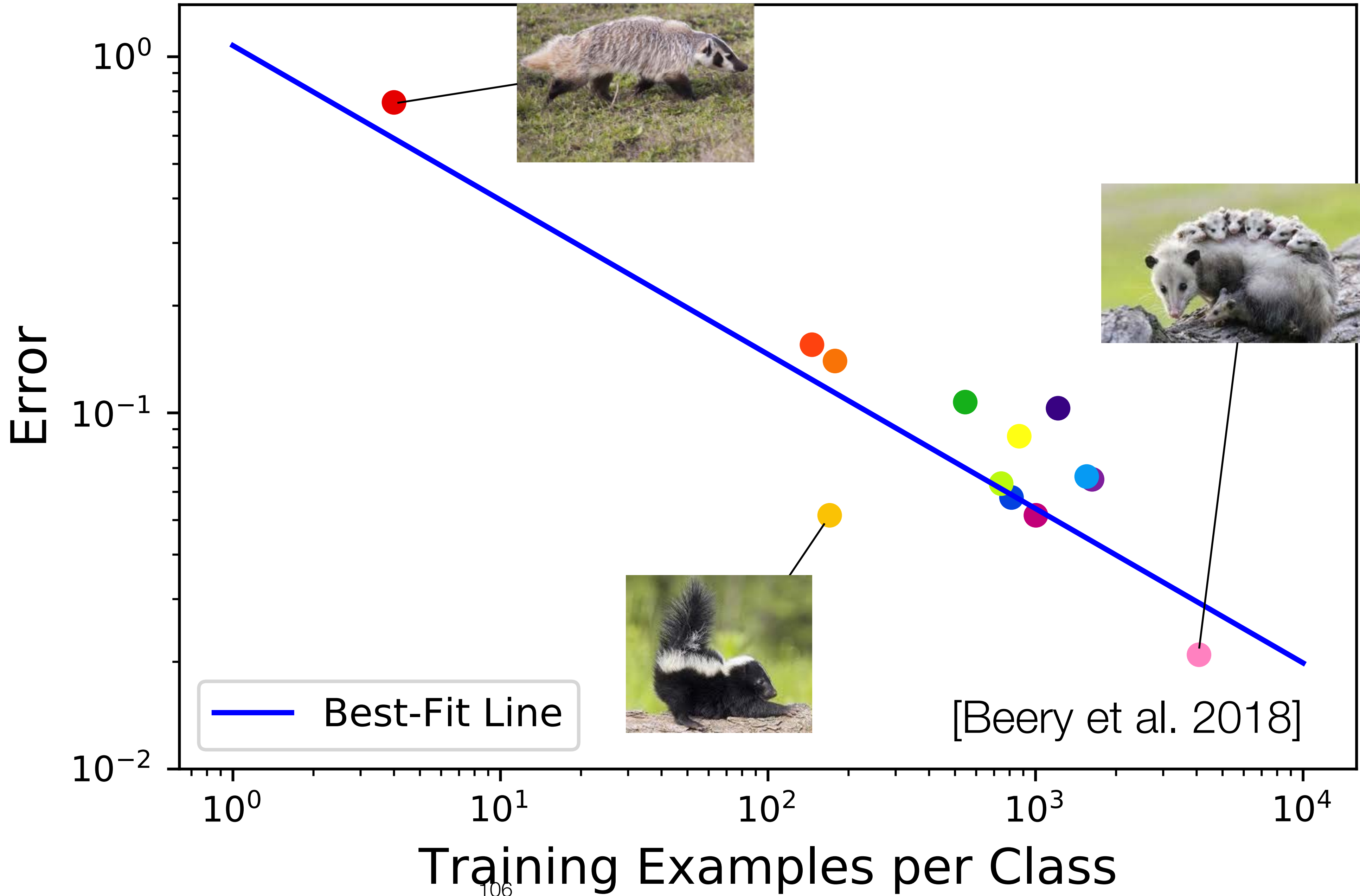
$y = \text{reward, outcome}$

$$p(y = 1|x, c) = p(y = 1|x)$$

$$p(y = 1|x, c = 0) = p(y = 1|x, c = 1)$$

- fox
- badger
- deer
- rodent
- skunk
- bird
- dog
- squirrel
- car
- cat
- bobcat
- coyote
- rabbit
- raccoon
- opossum

Frame Level Bbox Performance: Cis-Locations



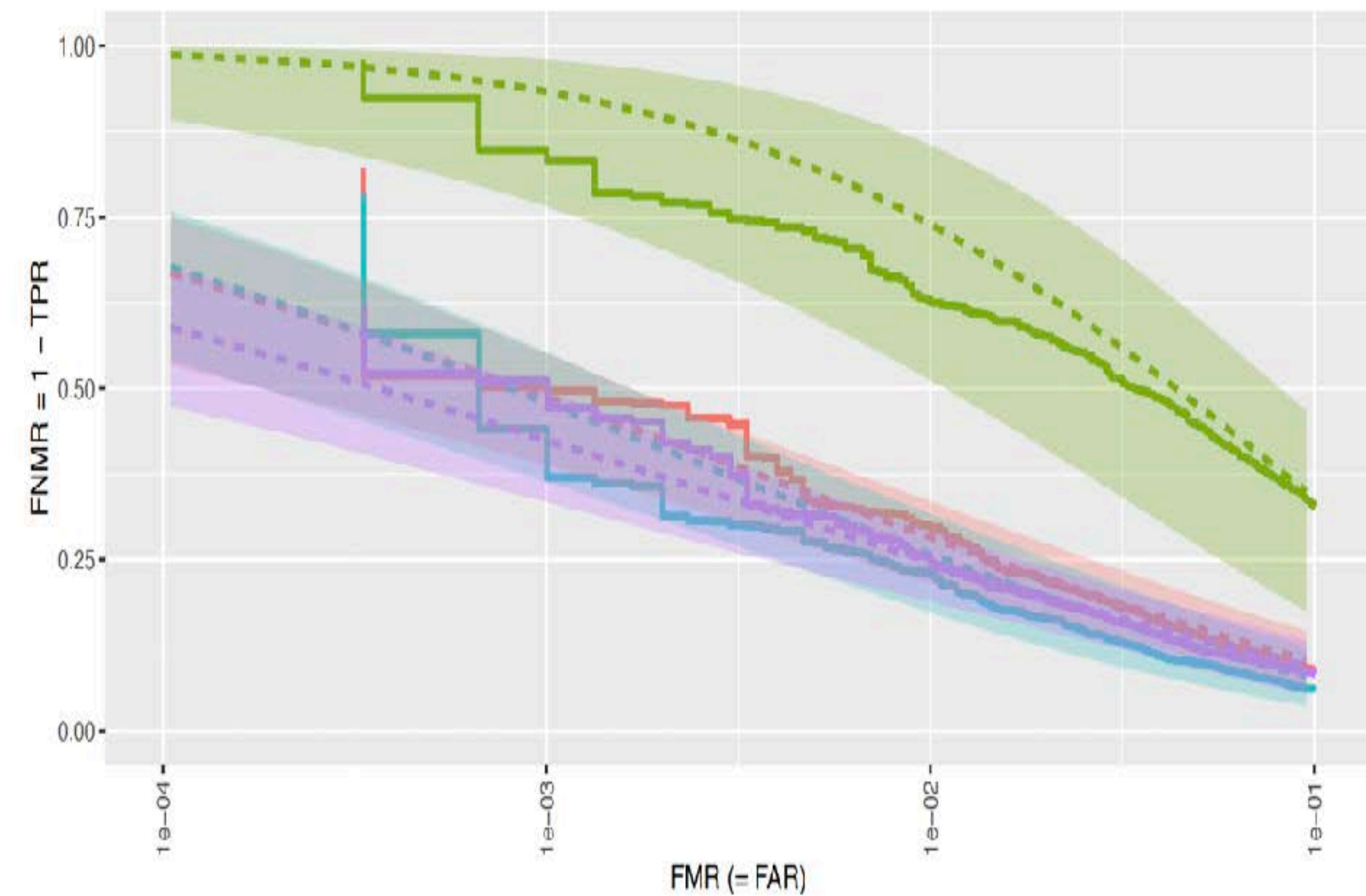
Unsupervised and Semi-supervised Bias Benchmarking in Face Recognition

Alexandra Chouldechova^{†*} Siqi Deng[†] Yongxin Wang
 Wei Xia^{*} Pietro Perona

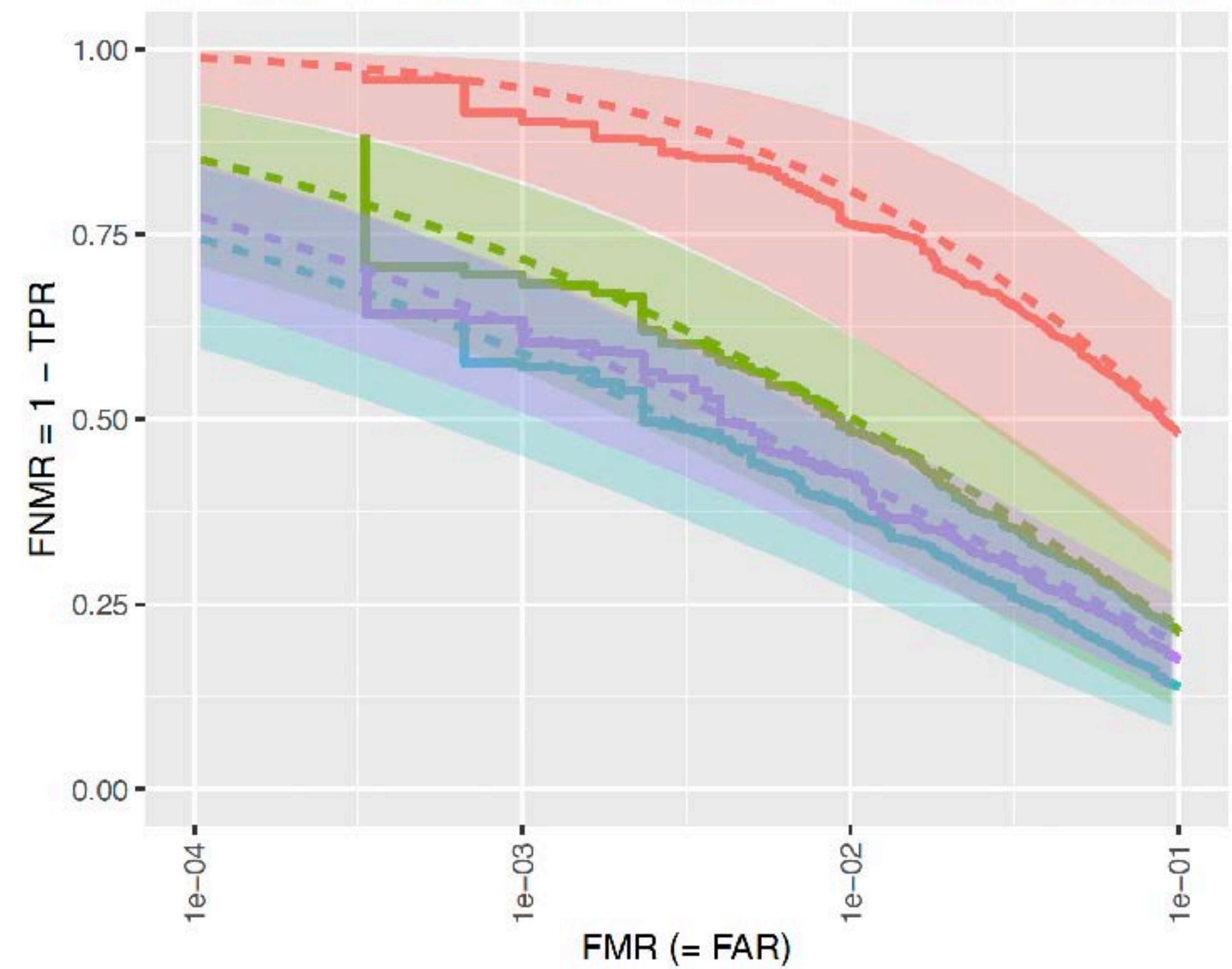
AWS AI Labs

ECCV'19

EA-RFW: SPE-FR TP-T FNMR vs FAR, NL = 0



AA_fresnet18_subarcface-RFW: SPE-FR TP-T FNMR vs:

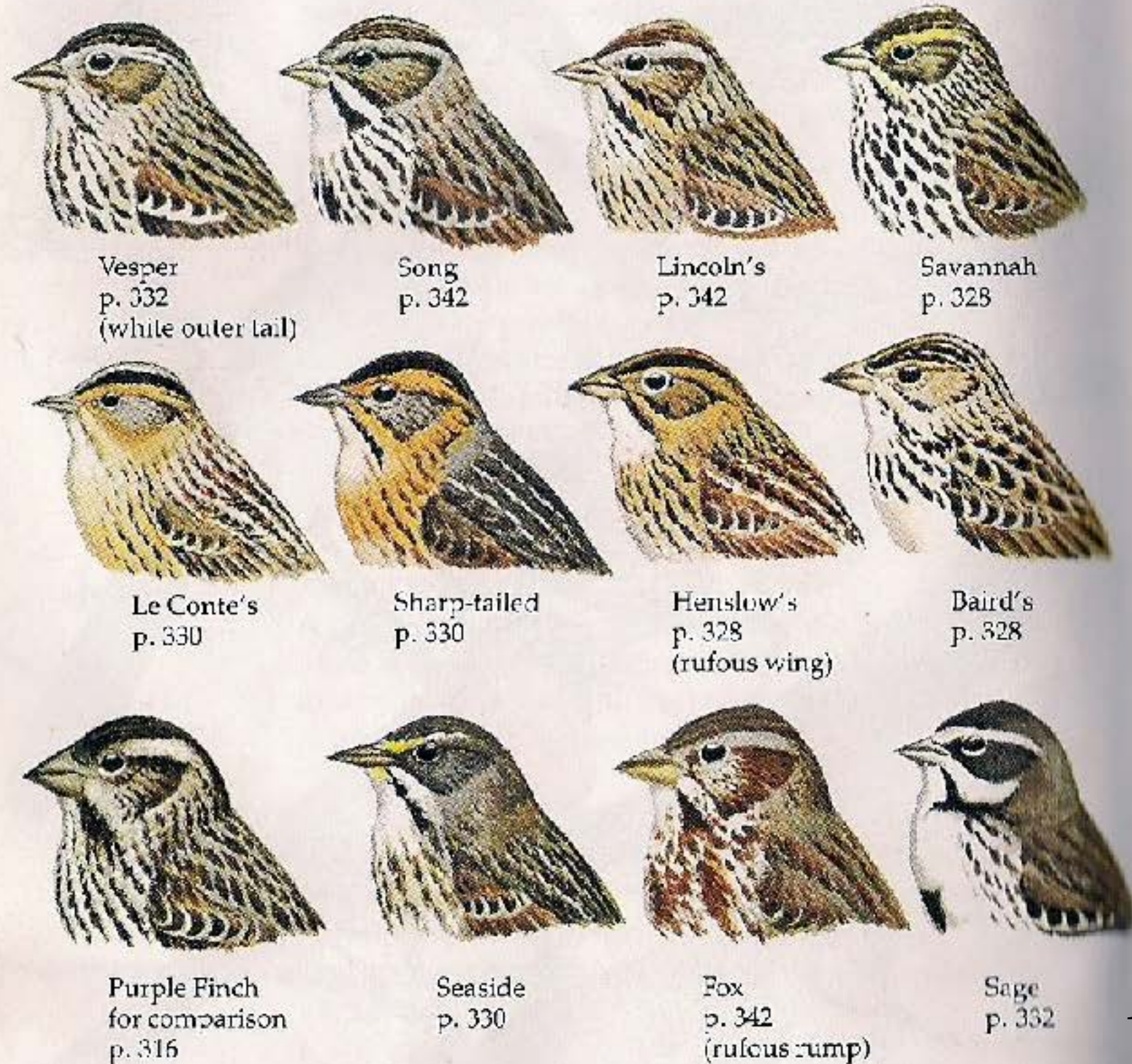


— African — Asian — Caucasian — Indian — True - - - SPE-f

• SPARROWS are small brown-bodied birds with streaked backs and short conical beaks. Their food, mostly seeds except during the nesting season, is obtained on or near the ground. When not nesting, most are seen in flocks. These may be identified by their song. The following species are shown by one head each; the striking Lark bunting, the Snow bunting, the towhees, and the Olive Sparrow are omitted. Immatures of some species are much duller, especially those species with black or rufous on the head. Songs and chips of sparrows are often more easily distinguished than are their plumages. See pp. 328-345 for further details.

Sparrows

STREAKED BREASTS



UNSTREAKED BREASTS



Summary

- Responsible AI
- Bias and fairness
- Correlation vs causation
- Bias in face analysis and recognition
- Bias in generative models
- Fixing bias
- Conclusions

You are in the top 0.05% for training, and talent....

Your responsibilities

- Know, and follow, the law
- Aim for big positive outcomes (people, the planet)
- Be informed. Think about both good and bad outcomes of your work
- Engage with society: teach, debate, get involved