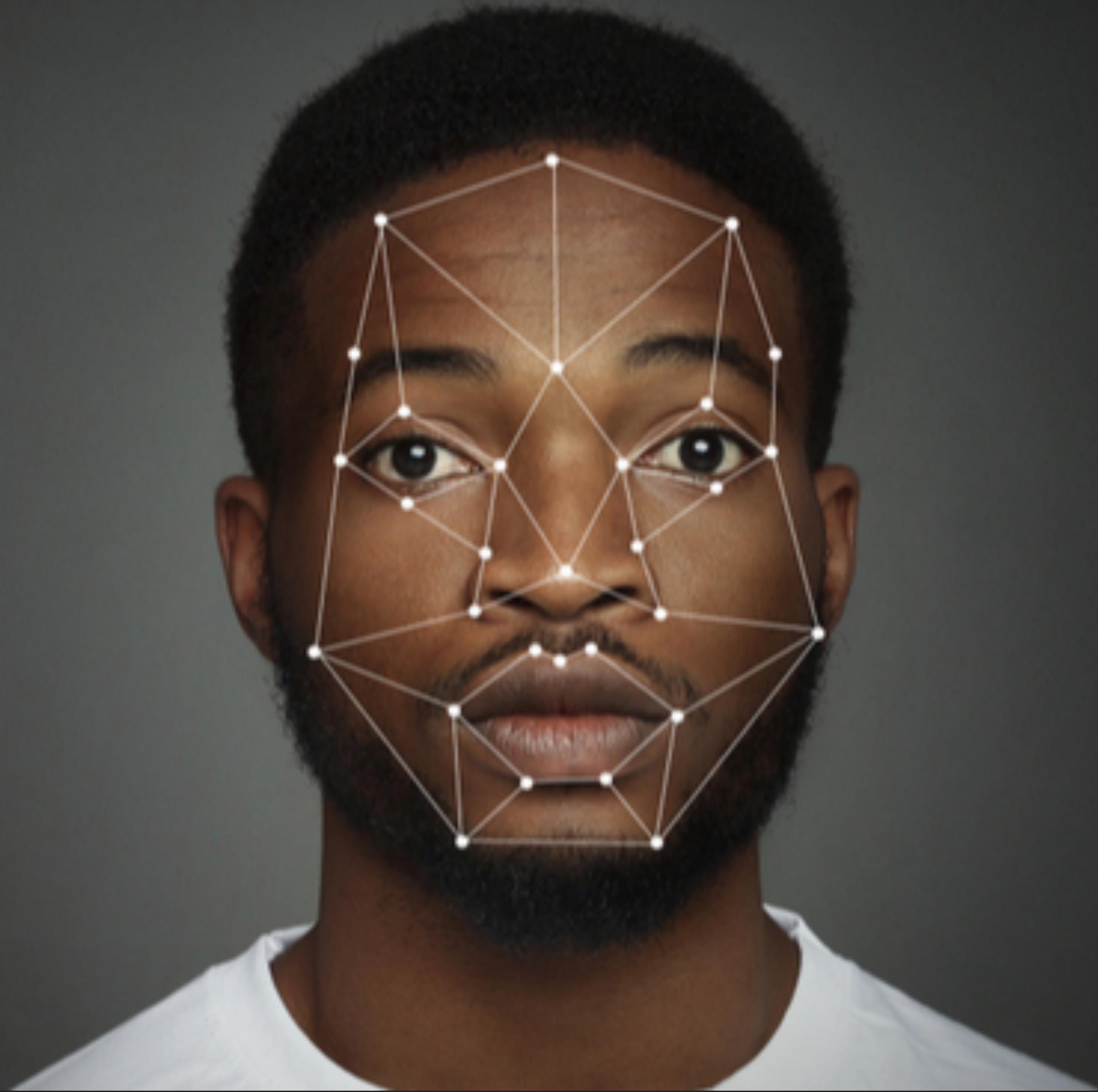


# **(Un)Ethical Computing**

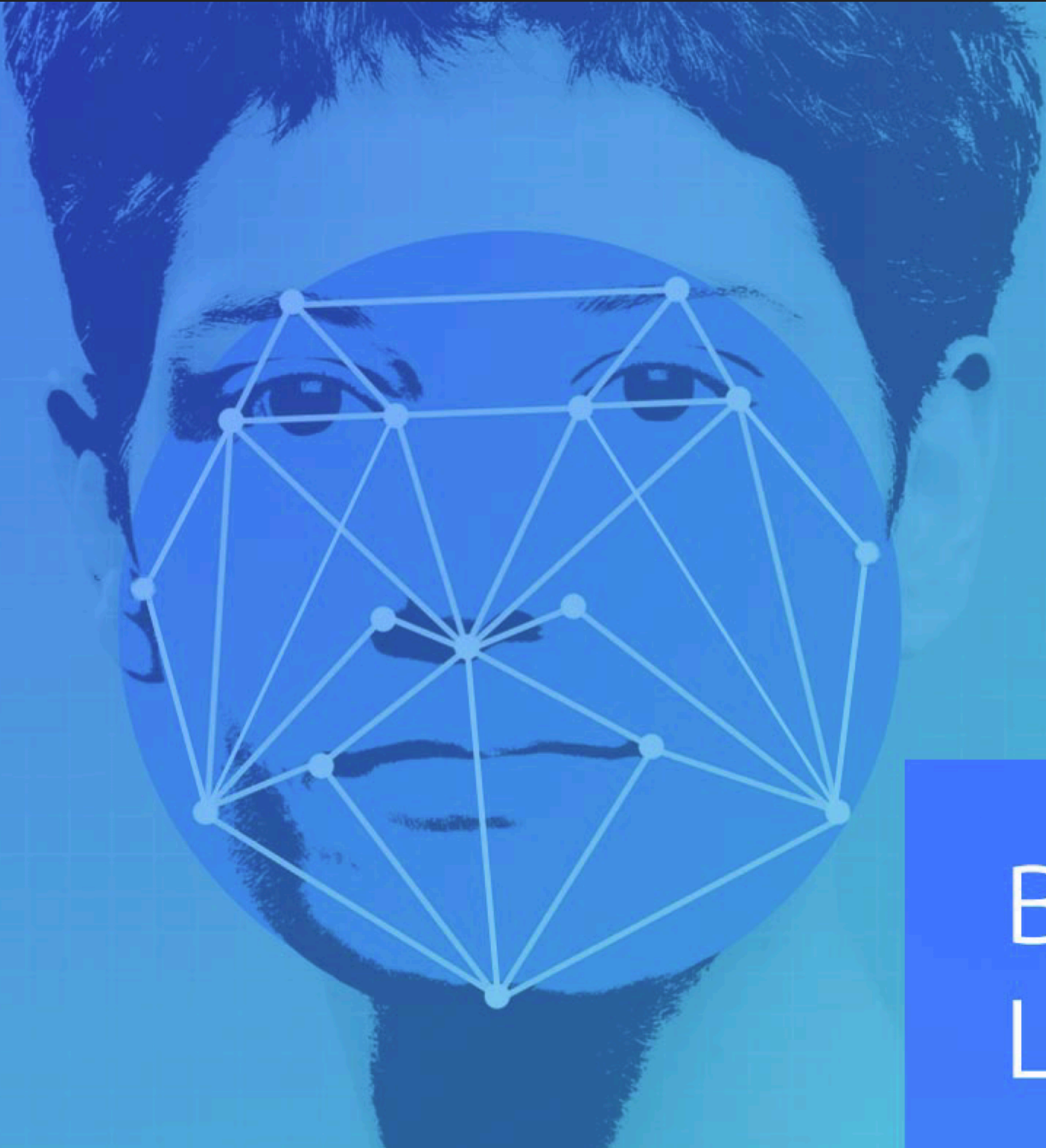
# face recognition



# predicting recidivism



# human resources



Biases in Machine  
Learning Algorithms

# financial lending



# non-consensual imagery

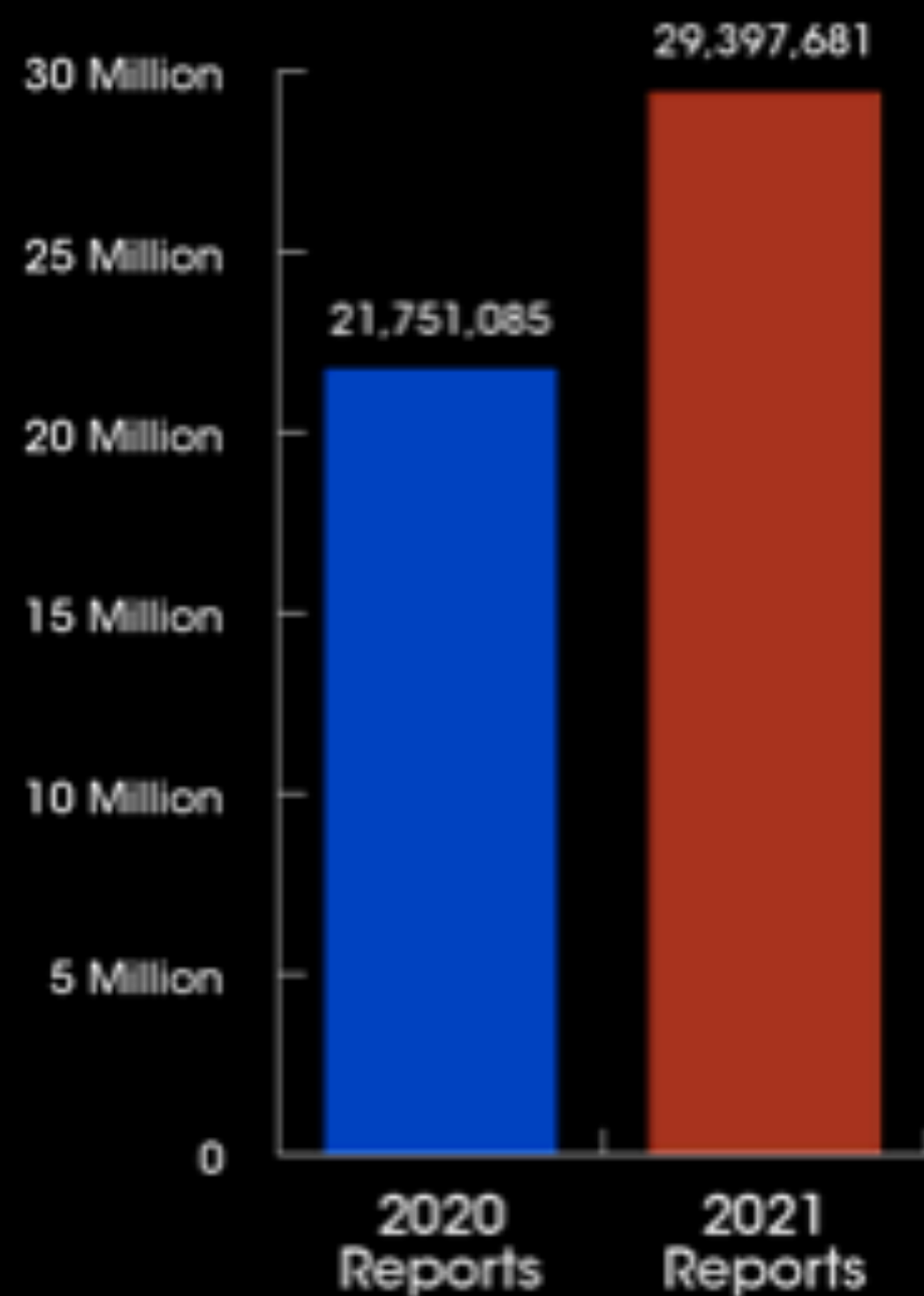


**addiction (body image, self harm, bullying,...)**



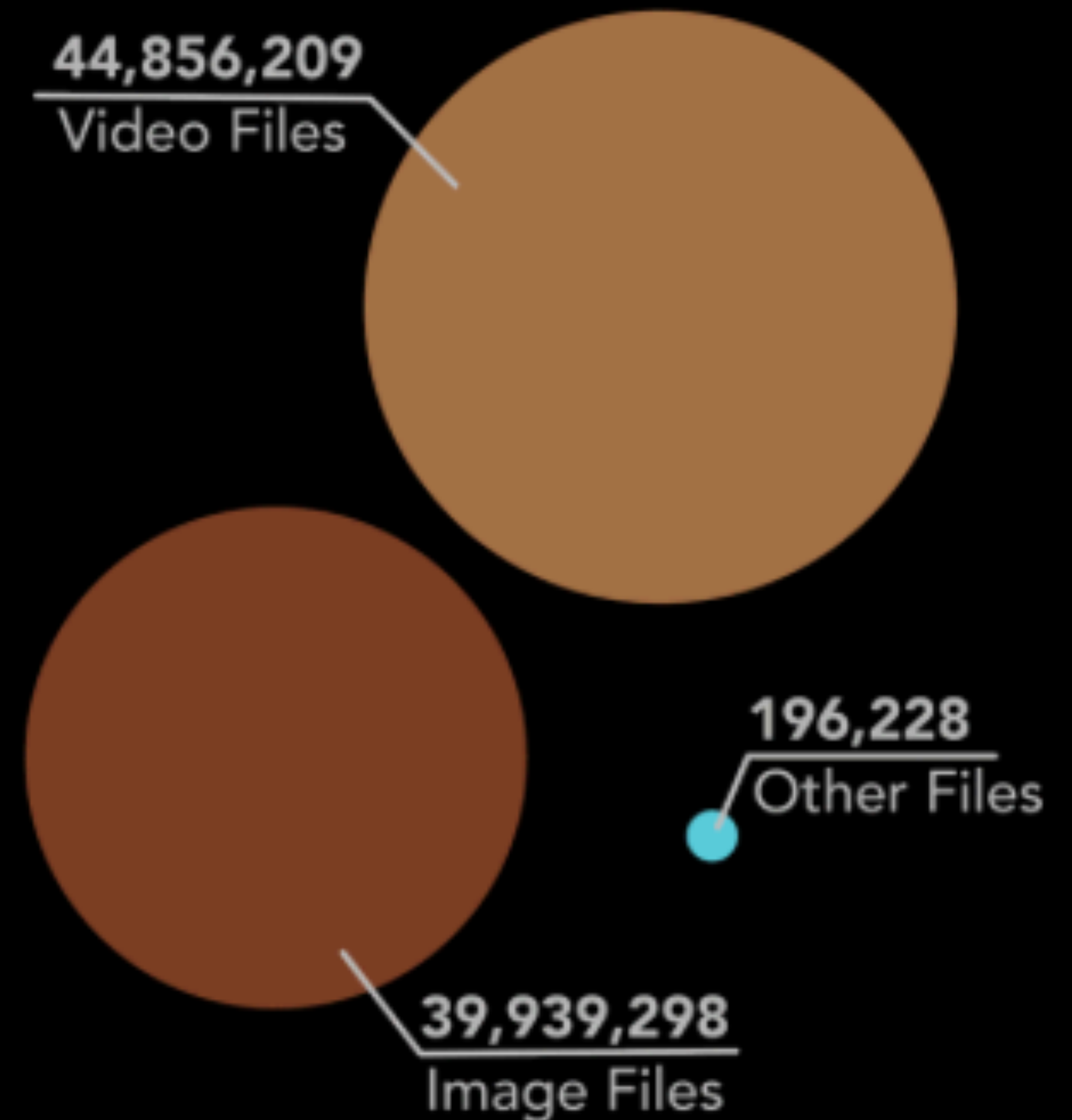
# child sexual abuse

## Total Reports



In 2021, reports to the CyberTipline increased by **35%** from 2020.

Reports to the CyberTipline included 85 million files



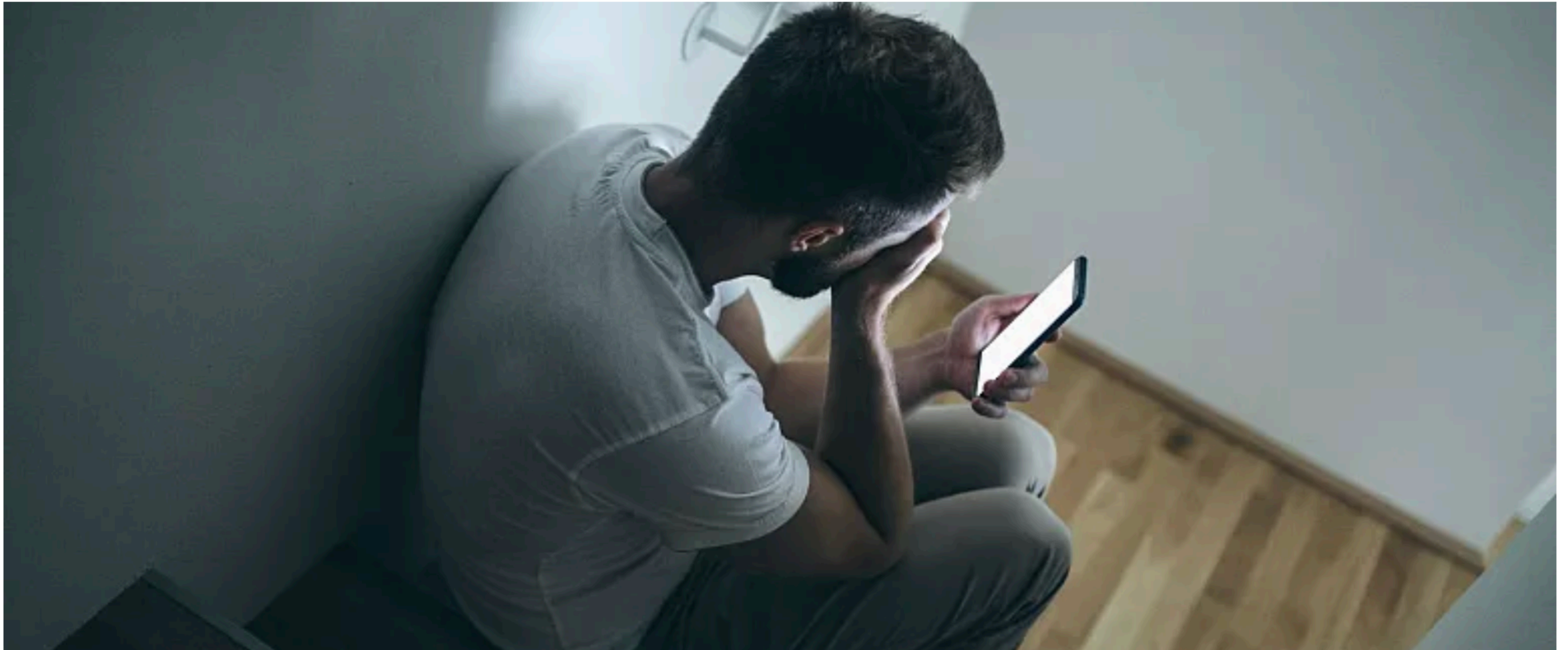


narcotics



## generative AI

**Man ends his life after an AI chatbot 'encouraged' him to sacrifice himself to stop climate change**



# predicting recidivism



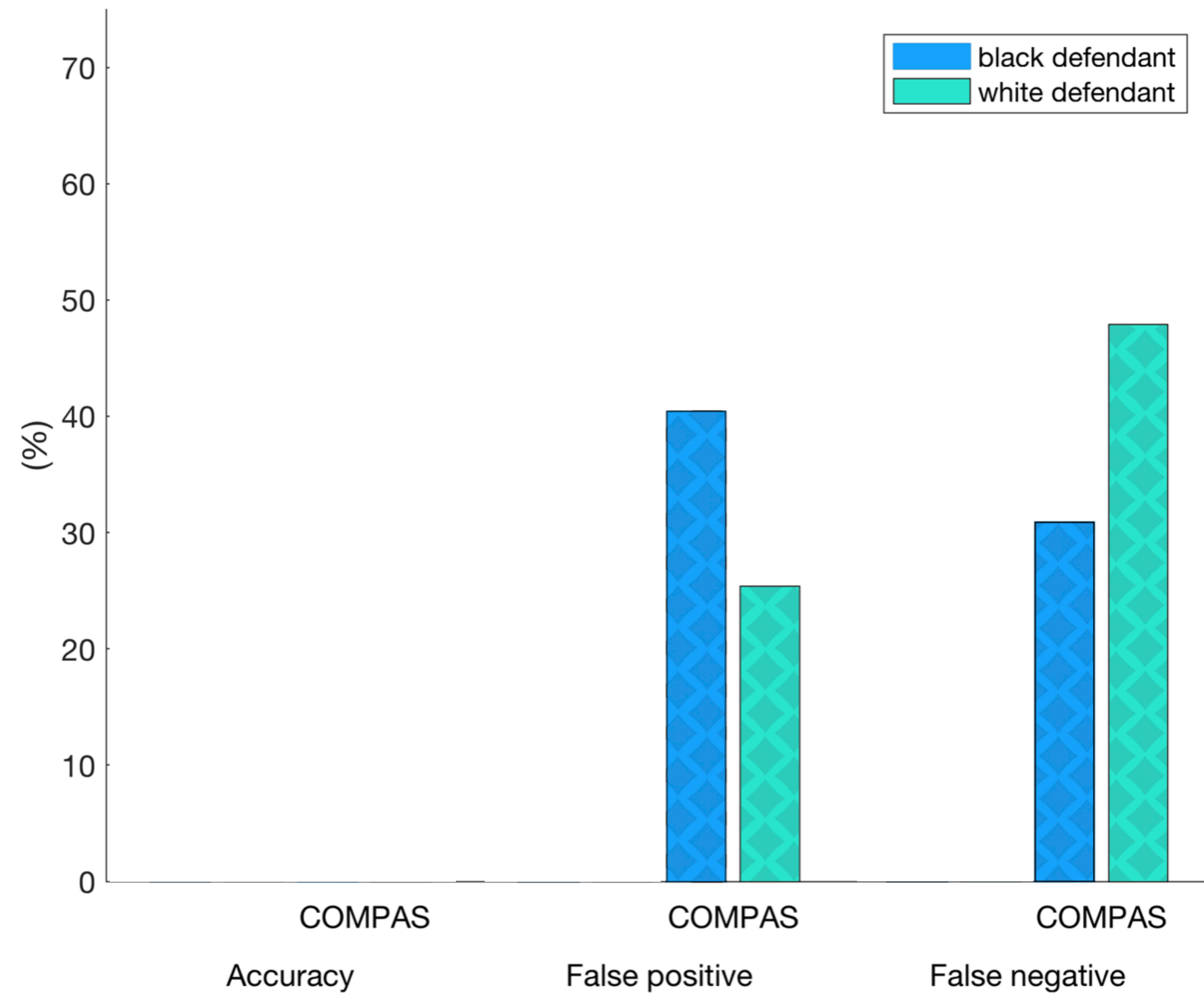


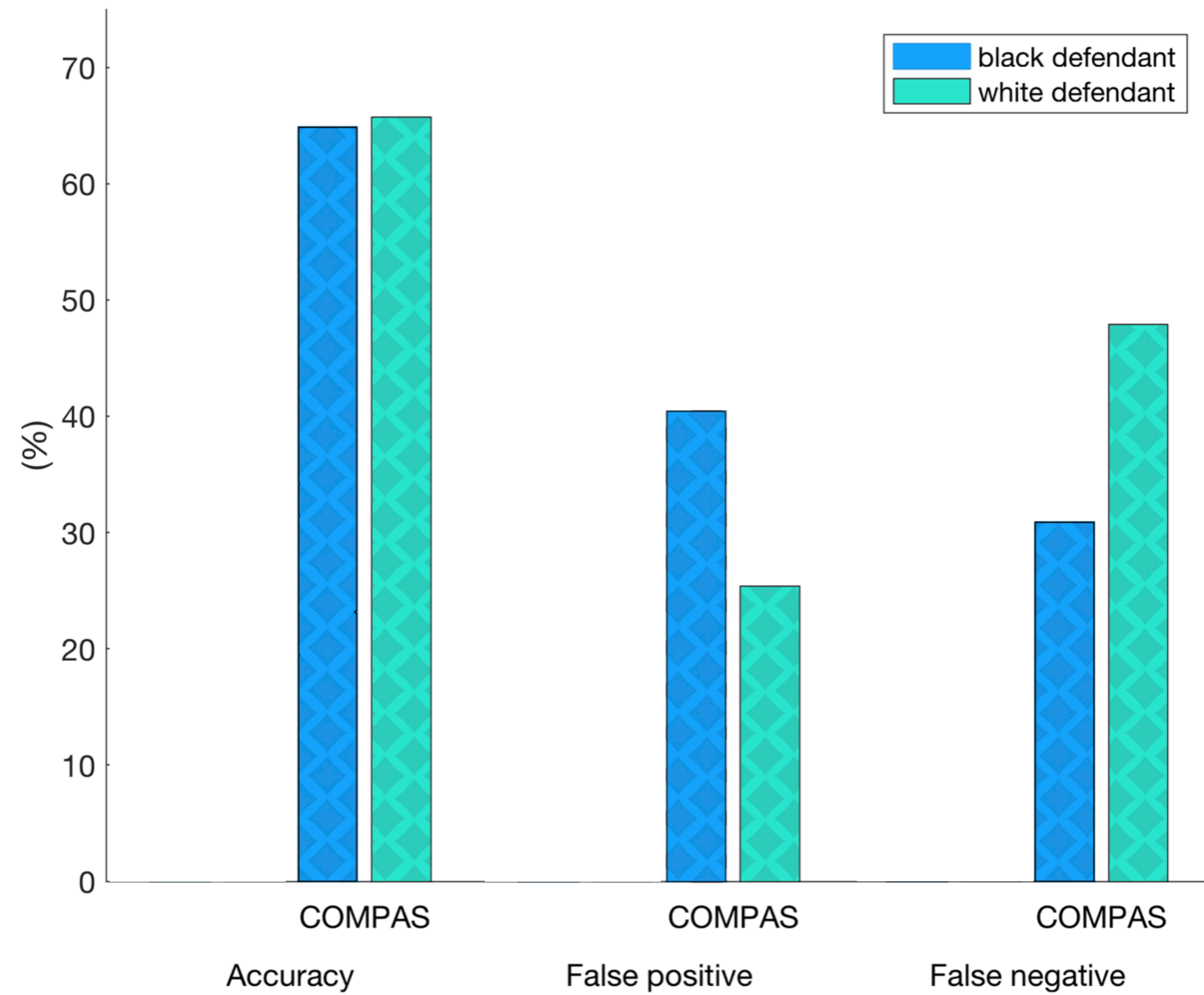
# Machine Bias

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





Your current accuracy is: 63%

The defendant is a **male** aged **43**. They have been charged with: **Disorderly Intoxication**. This crime is classified as a **misdemeanor**. They have been convicted of **2 prior crimes**. They have **0 juvenile felony charges** and **0 juvenile misdemeanor charges** on their record. Do you think this person will commit another crime within 2 years?

Yes

No

*Disorderly Intoxication: When an intoxicated person endangers the safety of another person or property, or causes a disturbance in public*



Your current accuracy is: 63%

The defendant is a **male** aged **43**. They have been charged with: **Disorderly Intoxication**. This crime is classified as a **misdemeanor**. They have been convicted of **2 prior crimes**. They have **0 juvenile felony charges** and **0 juvenile misdemeanor charges** on their record. Do you think this person will commit another crime within 2 years?

Yes

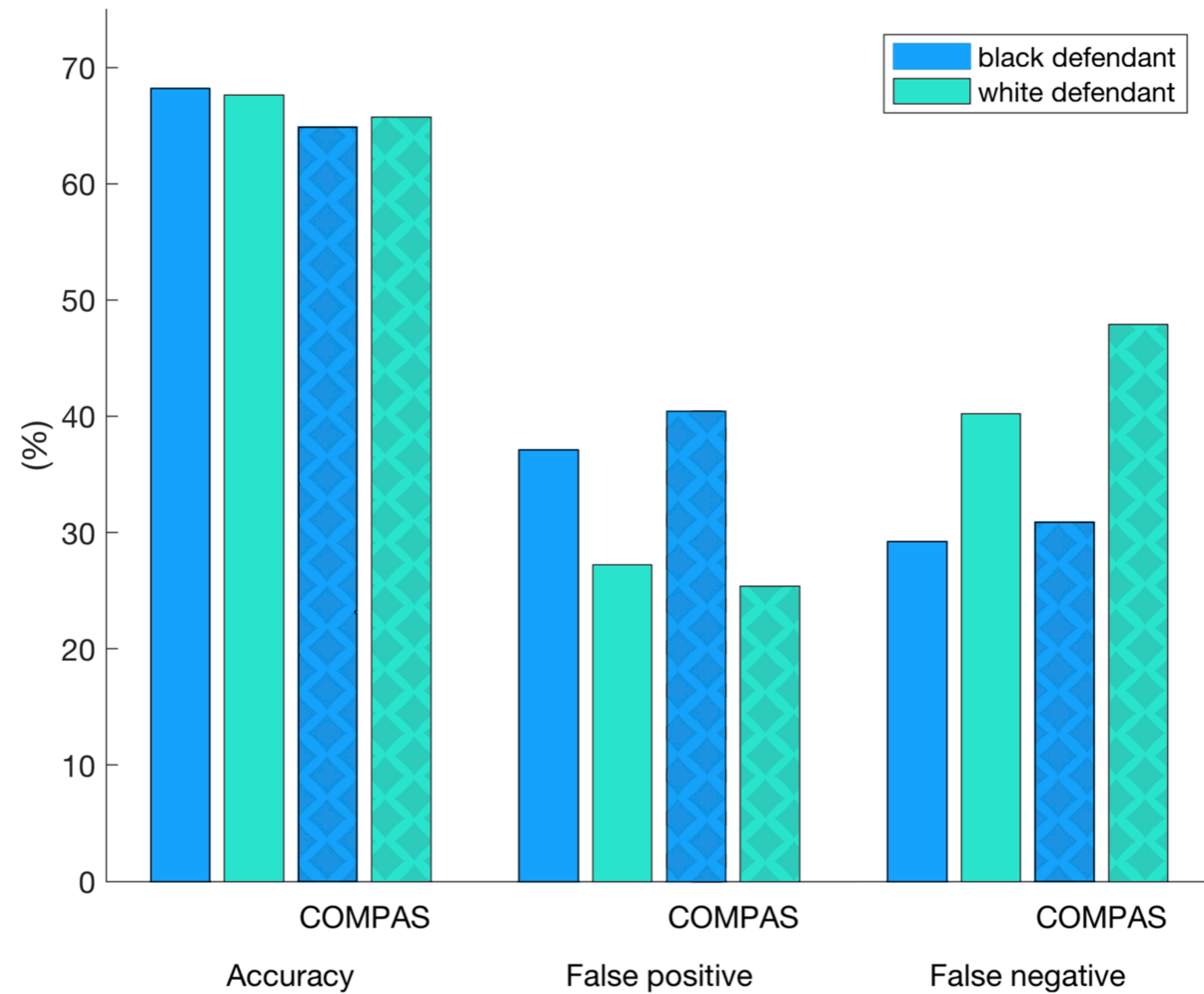
No

*Disorderly Intoxication: When an intoxicated person endangers the safety of another person or property, or causes a disturbance in public*

- N = 400
- 50 questions/participant
- feedback
- catch trials
- \$1 payment
- \$5 bonus for > 65%

>>





Your current accuracy is: 63%

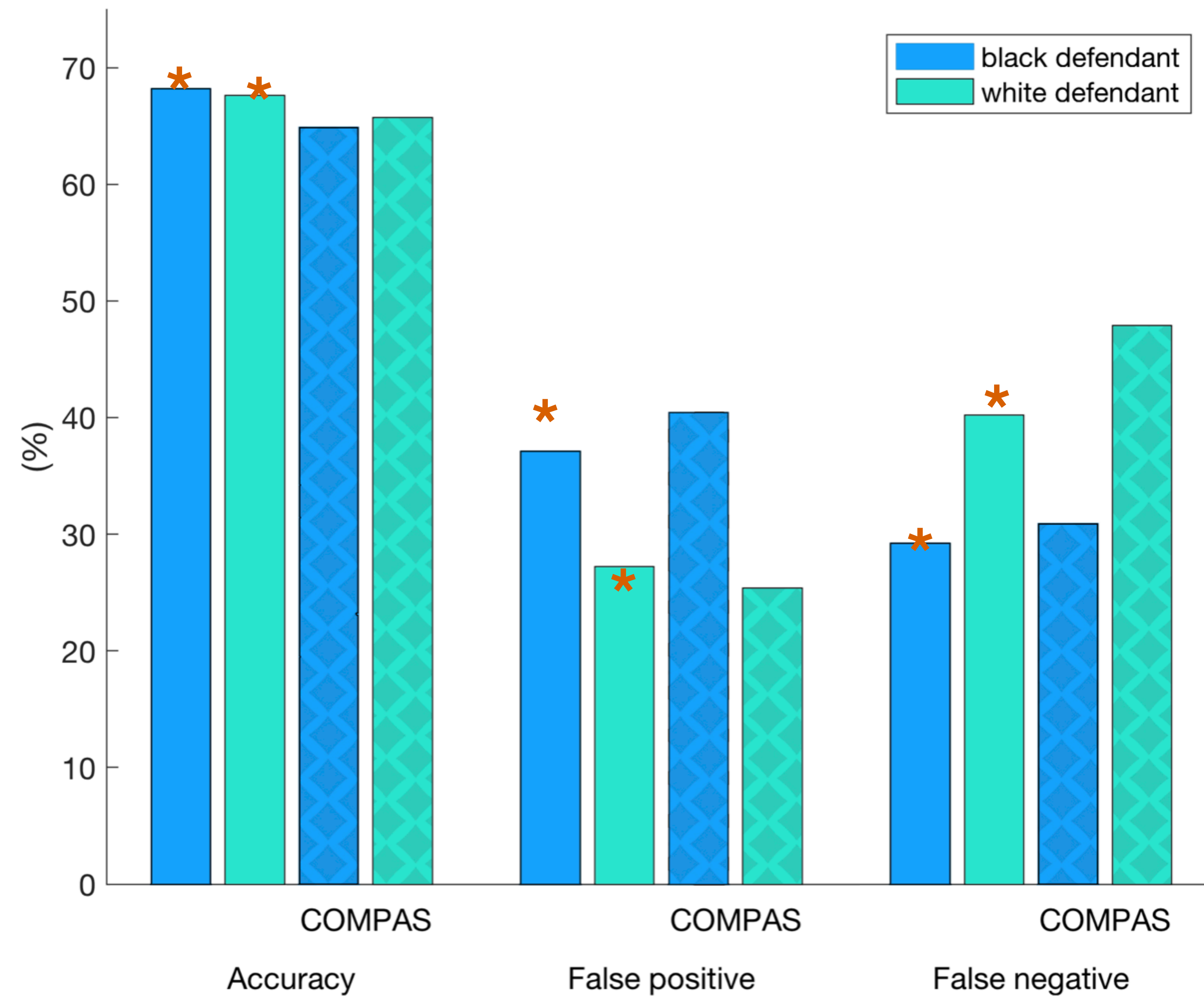
The defendant is a **white male** aged **43**. They have been charged with: **Disorderly Intoxication**. This crime is classified as a **misdemeanor**. They have been convicted of **2 prior crimes**. They have **0 juvenile felony charges** and **0 juvenile misdemeanor charges** on their record. Do you think this person will commit another crime within 2 years?

Yes

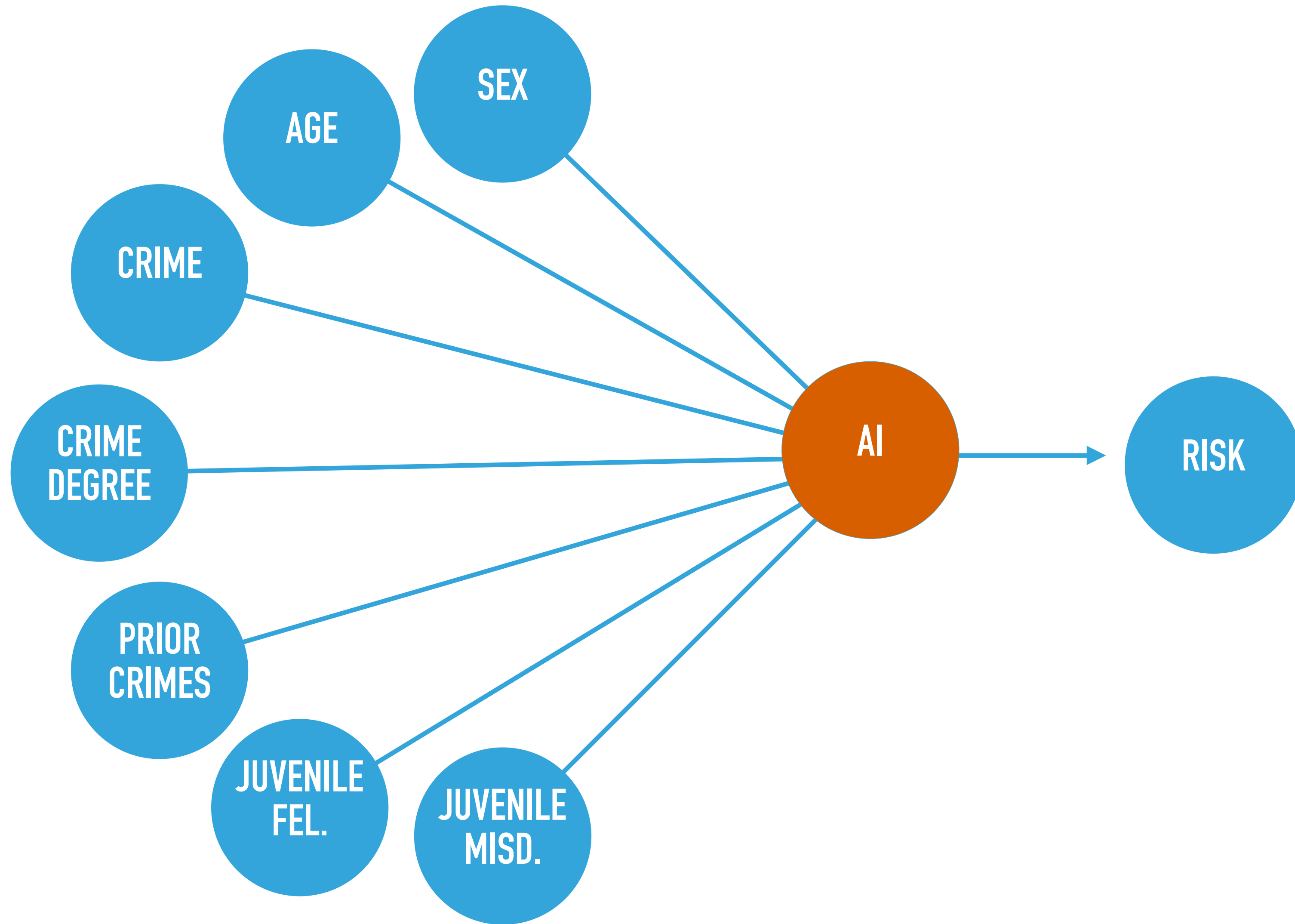
No

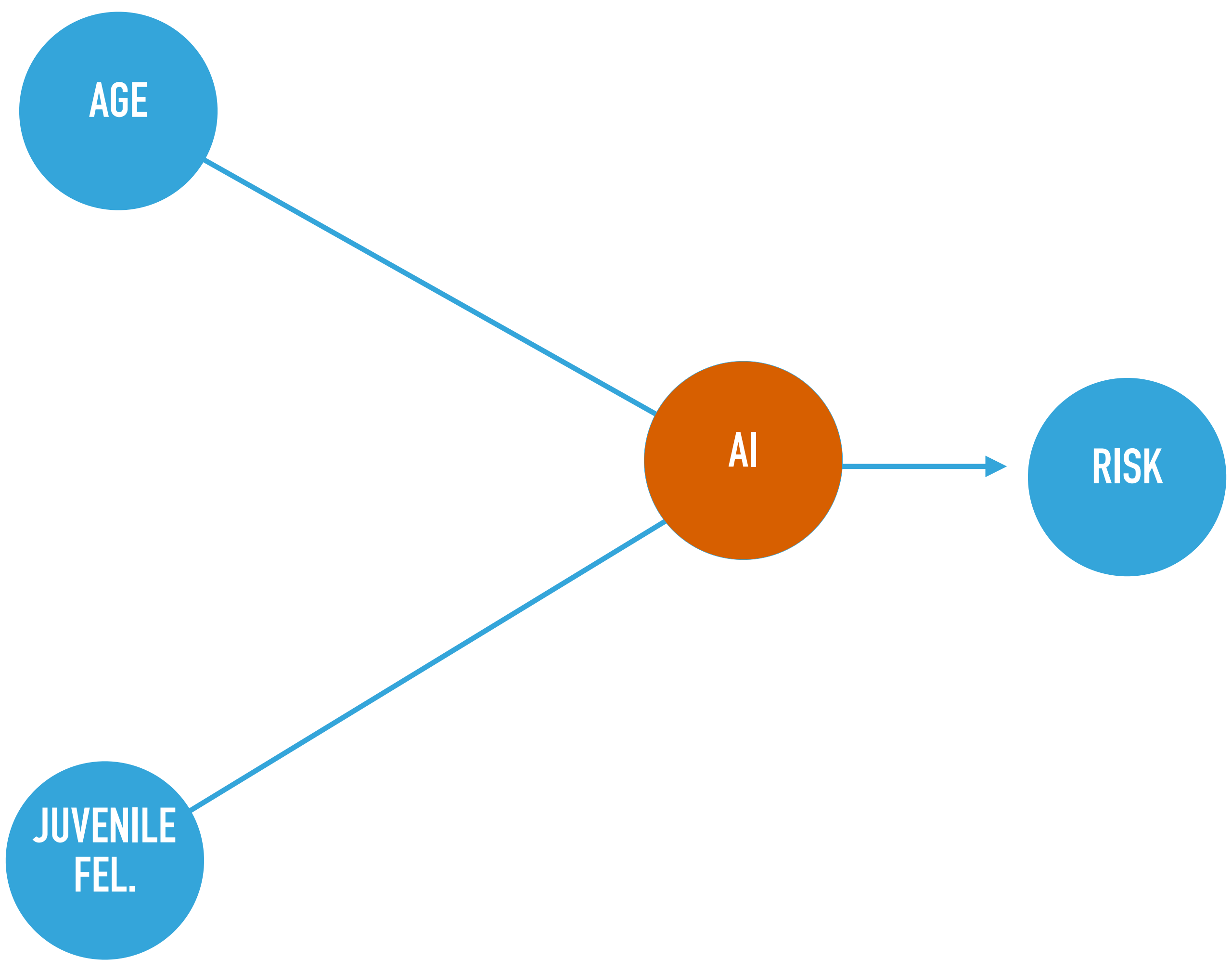
*Disorderly Intoxication: When an intoxicated person endangers the safety of another person or property, or causes a disturbance in public*

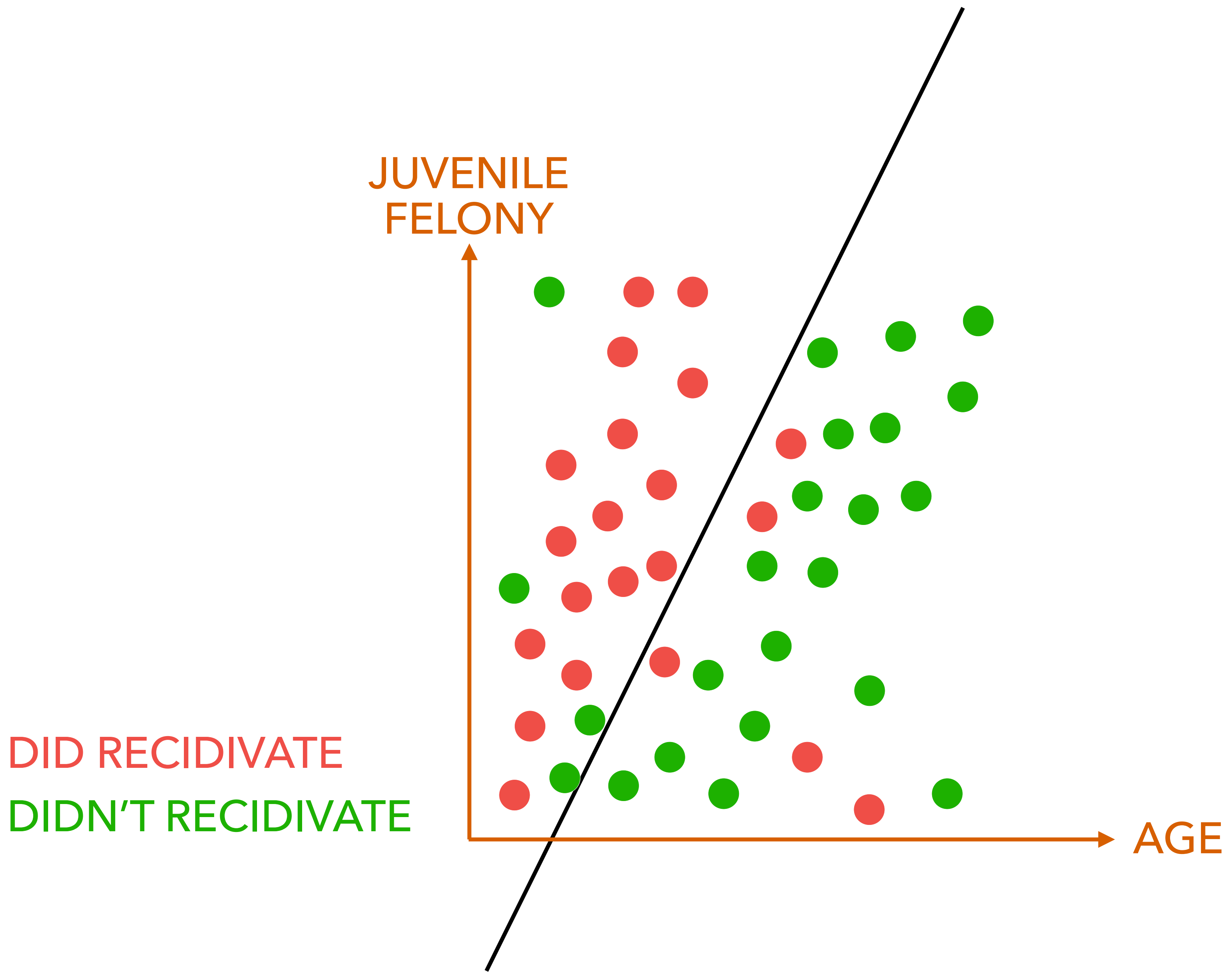


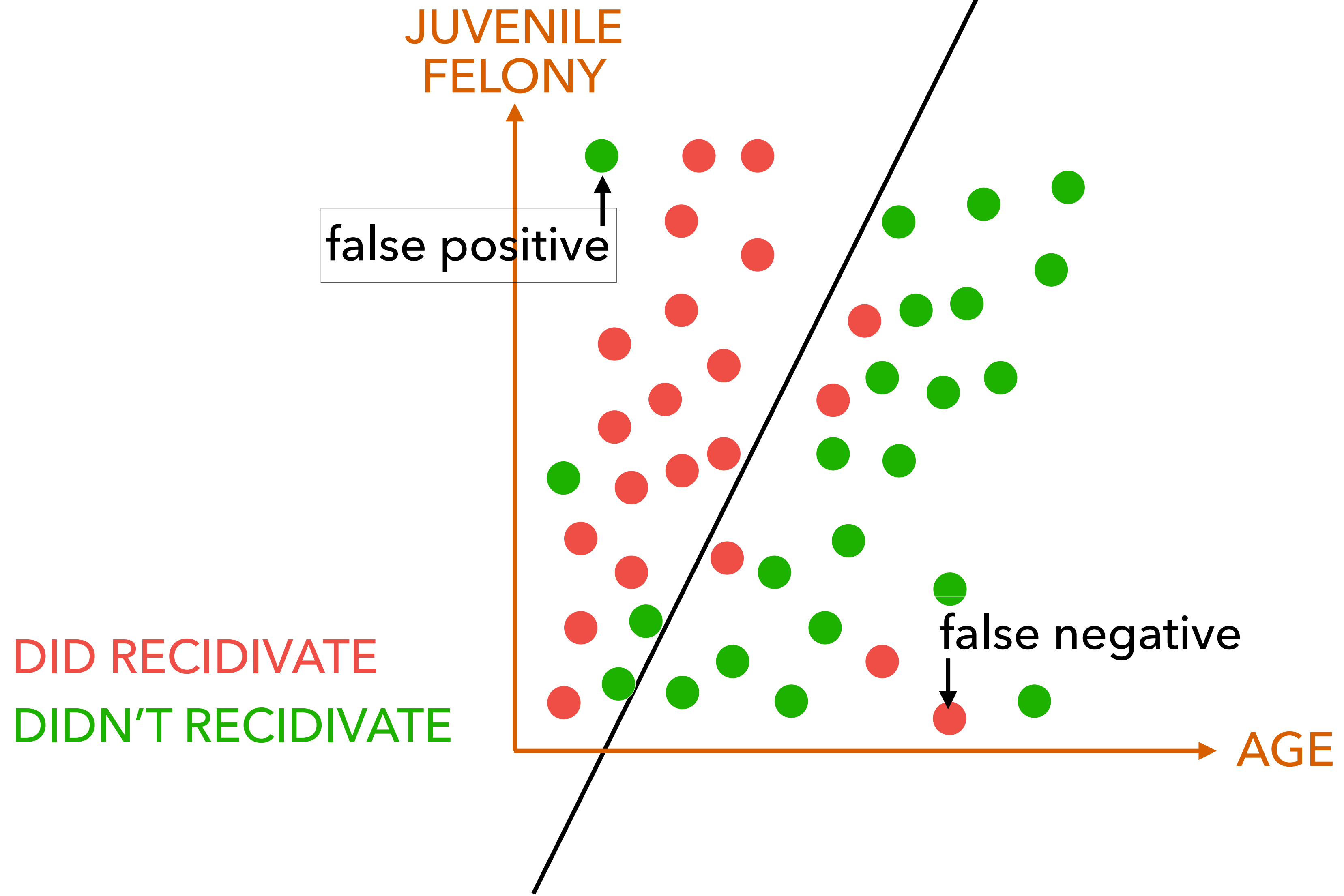


- 1. non-experts = AI?**
- 2. biased without race?**

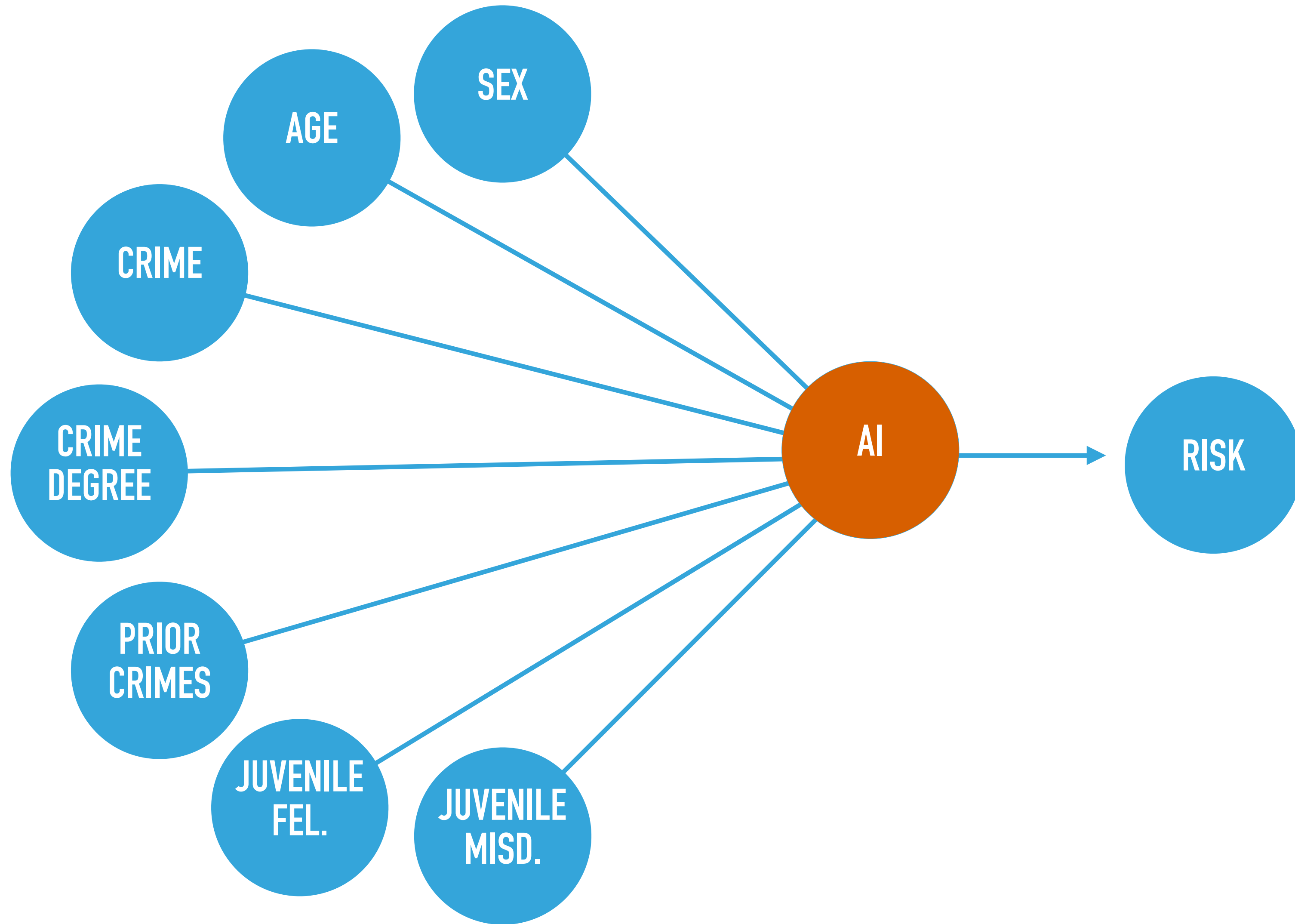


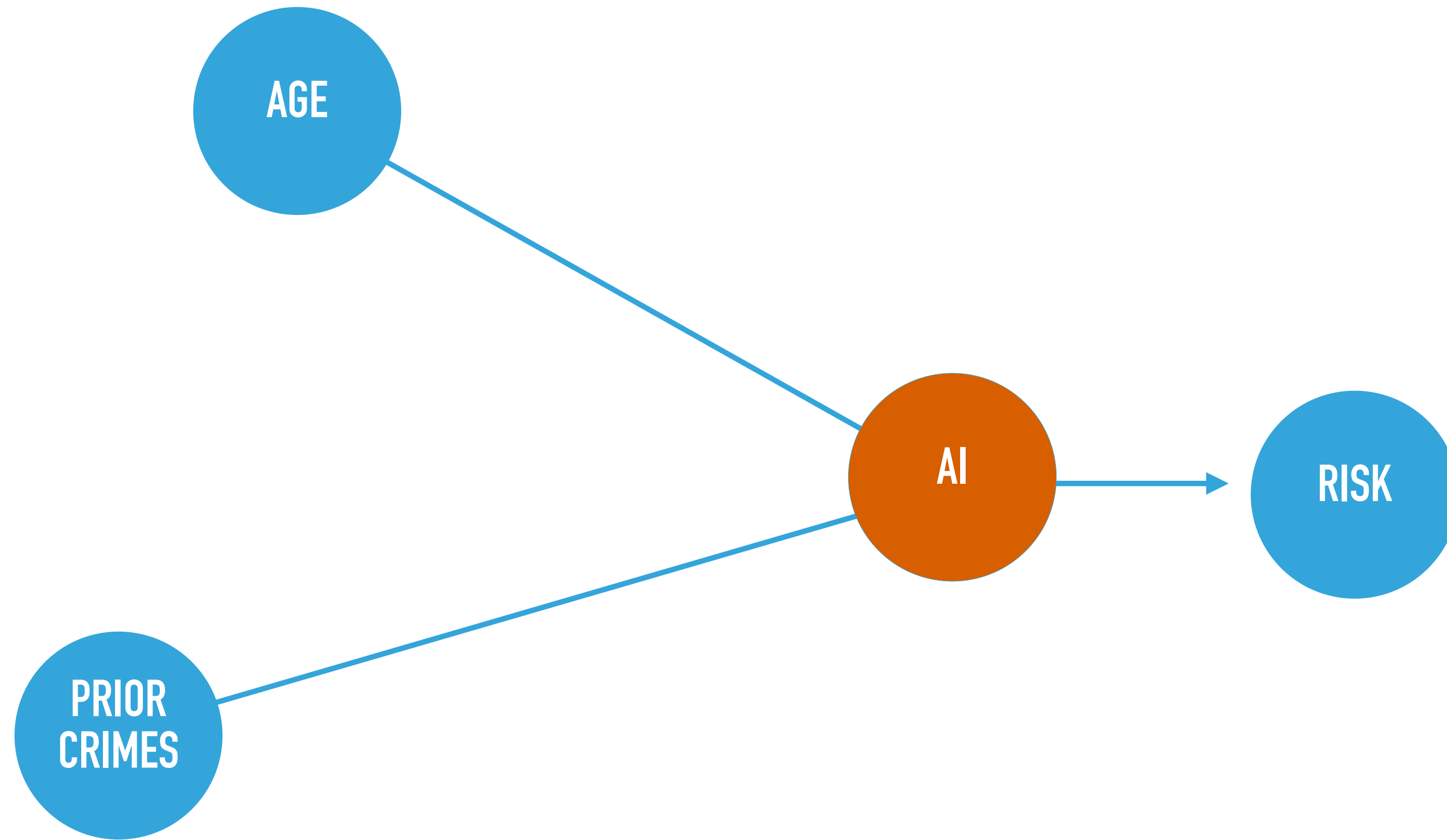


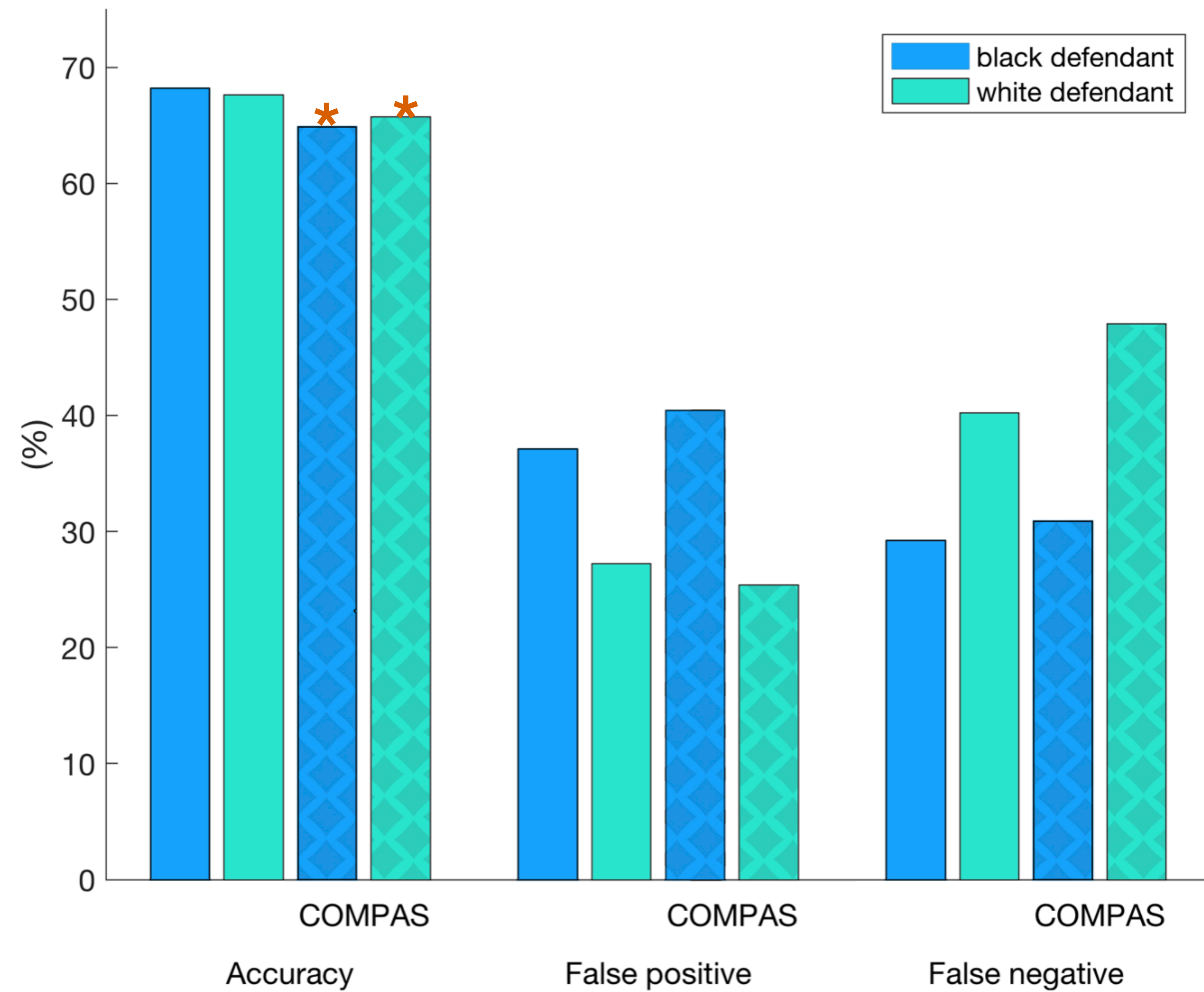


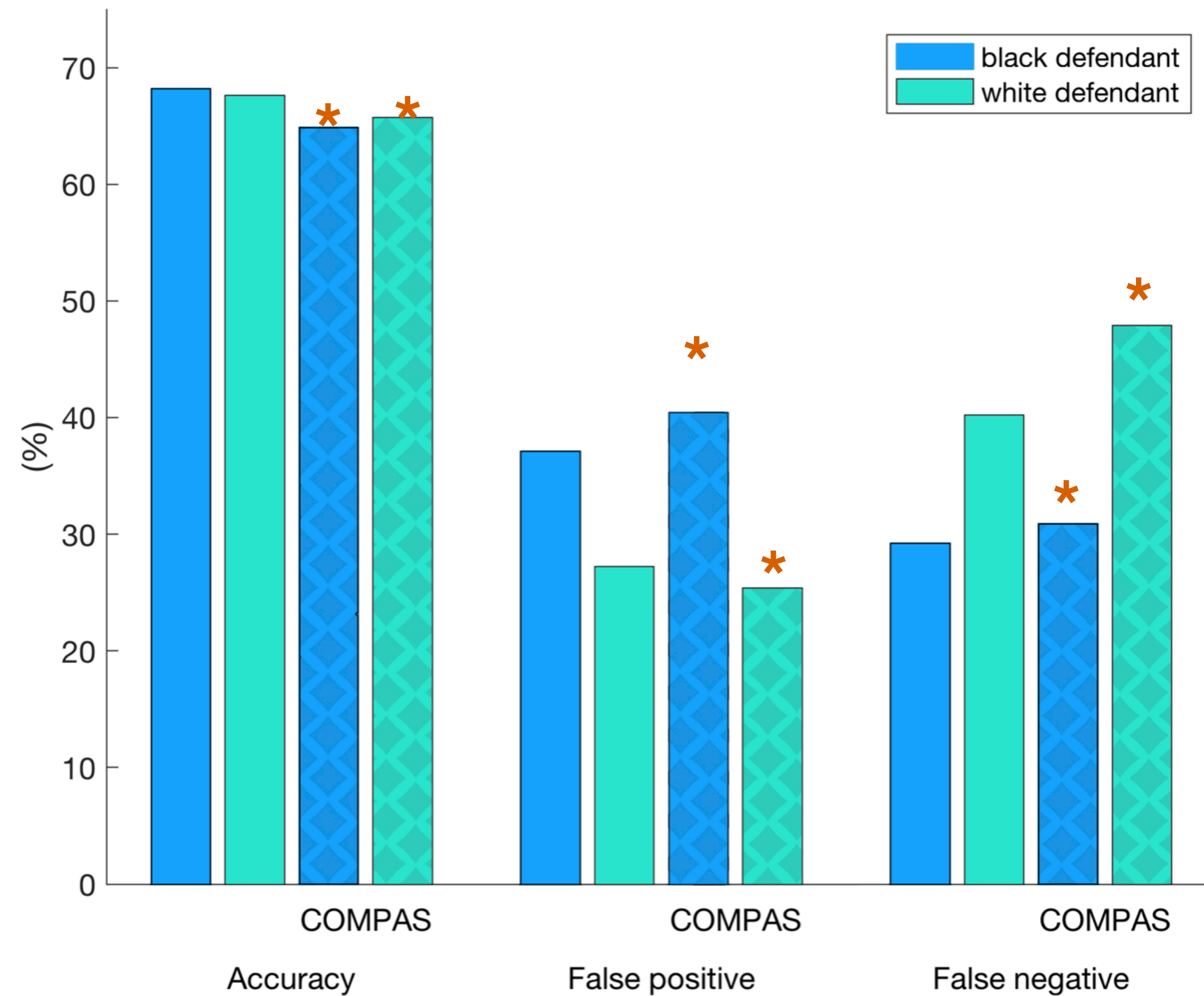






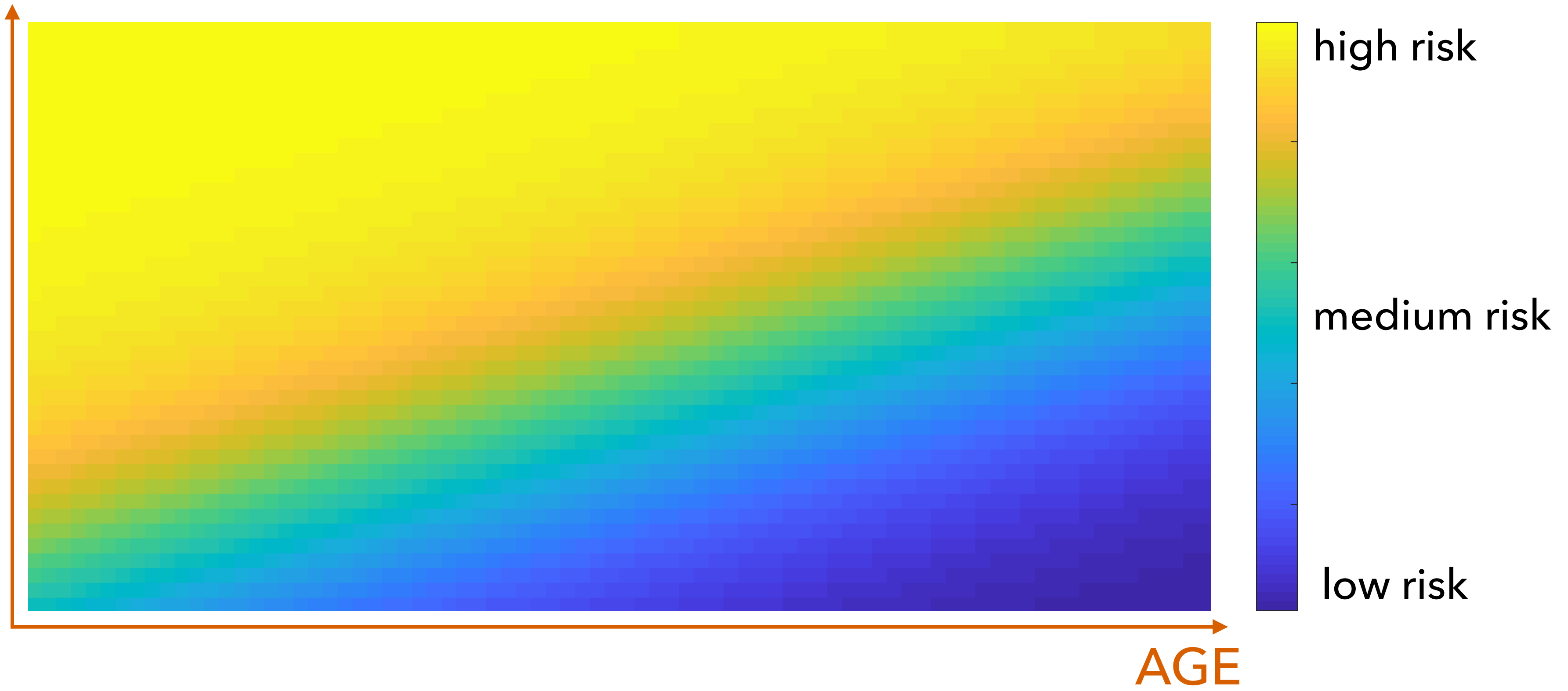






- 1. non-experts = AI?**
- 2. biased without race?**

PRIOR  
CRIMES



## **2. Prior Crimes is a proxy for Race**

**1. classification is easy (but not accurate)**



**so what?**

**Explainable AI/ML is important**

**The best way to repeat history is  
to train AI/ML on historical data**

**J. Dressel** and H. Farid. The Dangers of Risk Prediction in the Criminal Justice System. MIT Case Studies in Social and Ethical Responsibilities of Computing, 2021.



# non-consensual imagery





image-to-image translation  
(deep fakes)

start



appearance

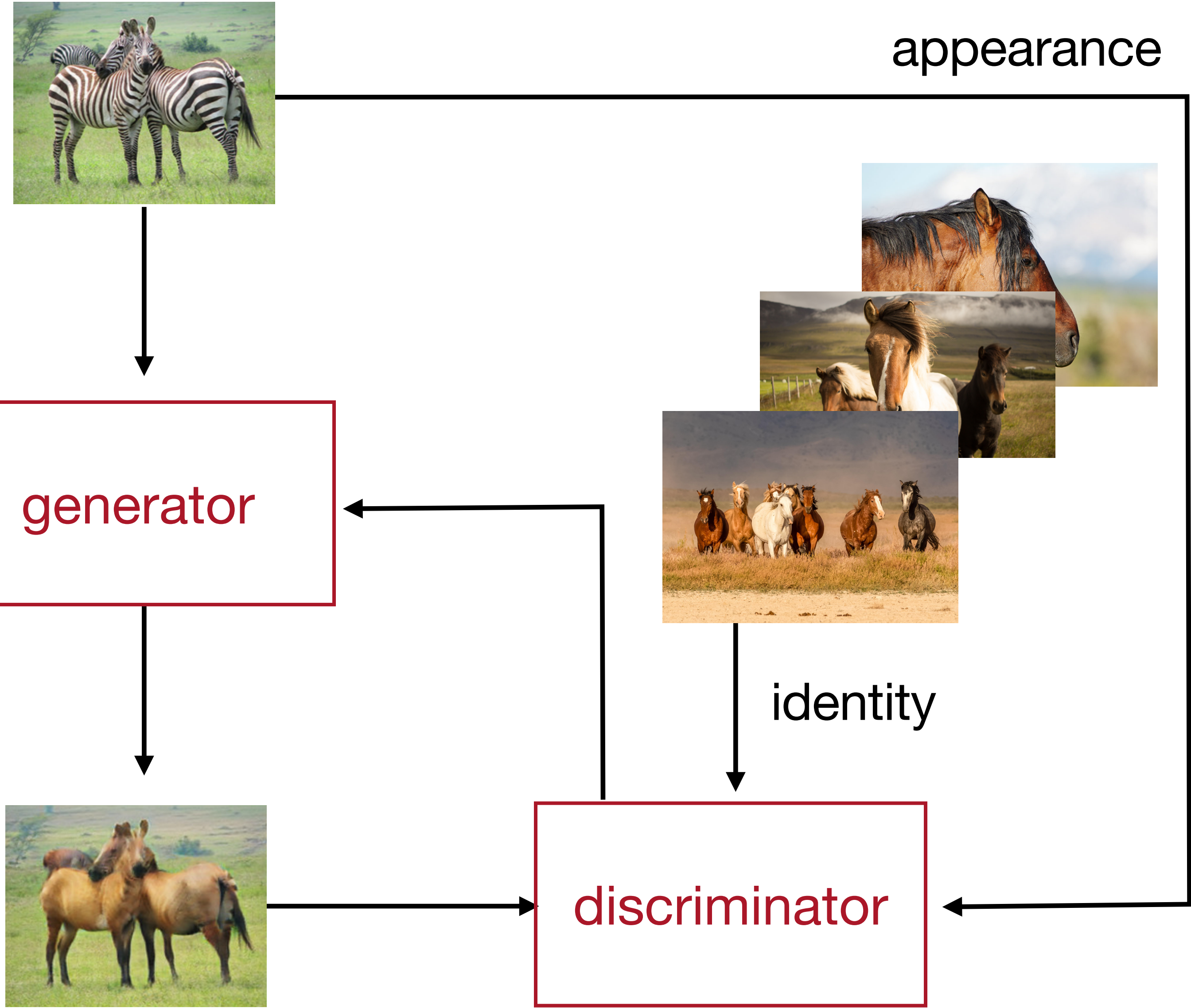


generator

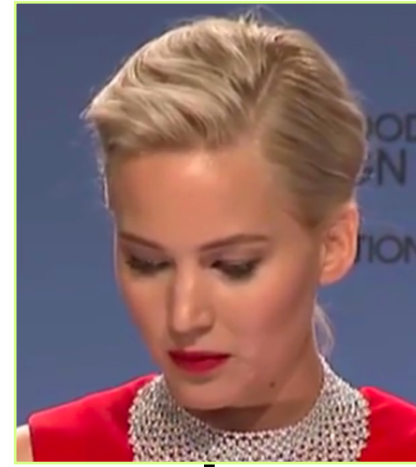


identity

discriminator

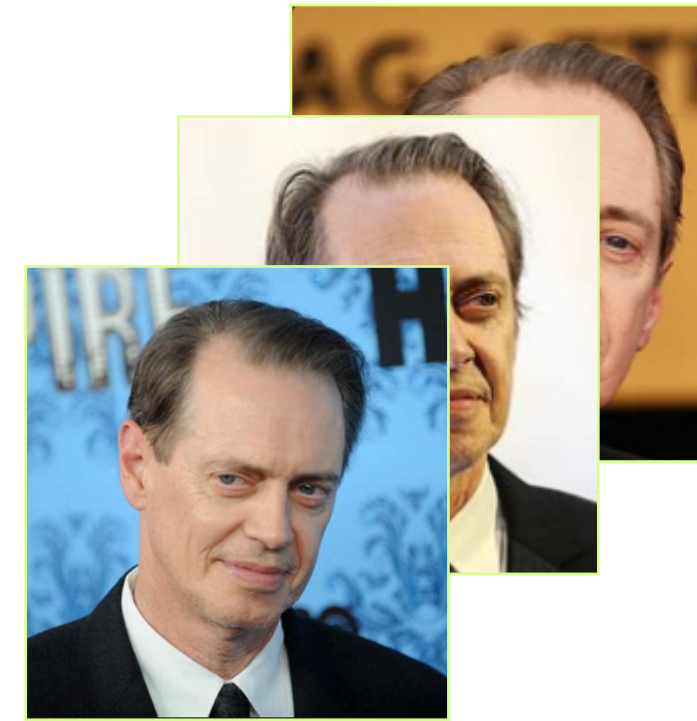


start



appearance

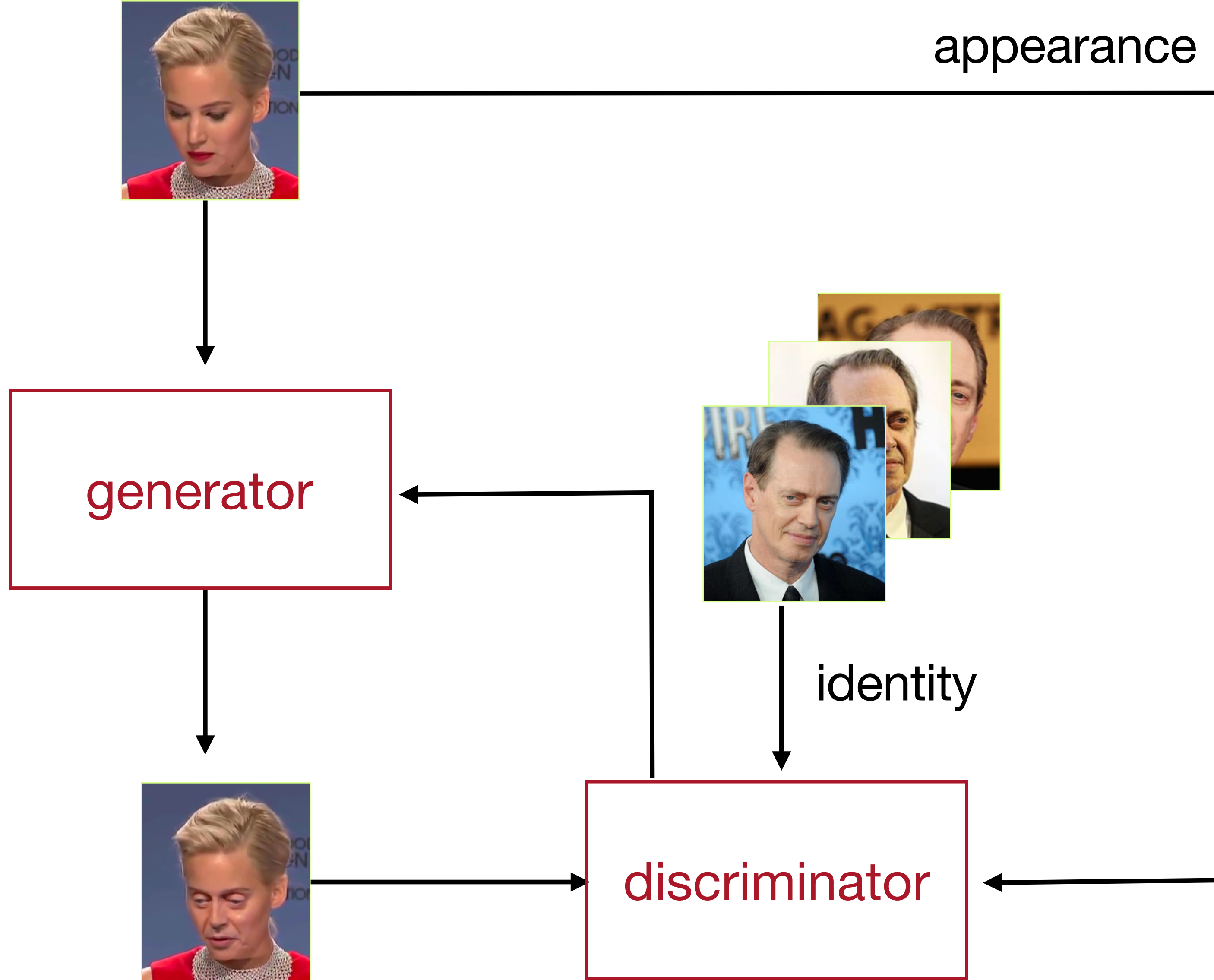
generator



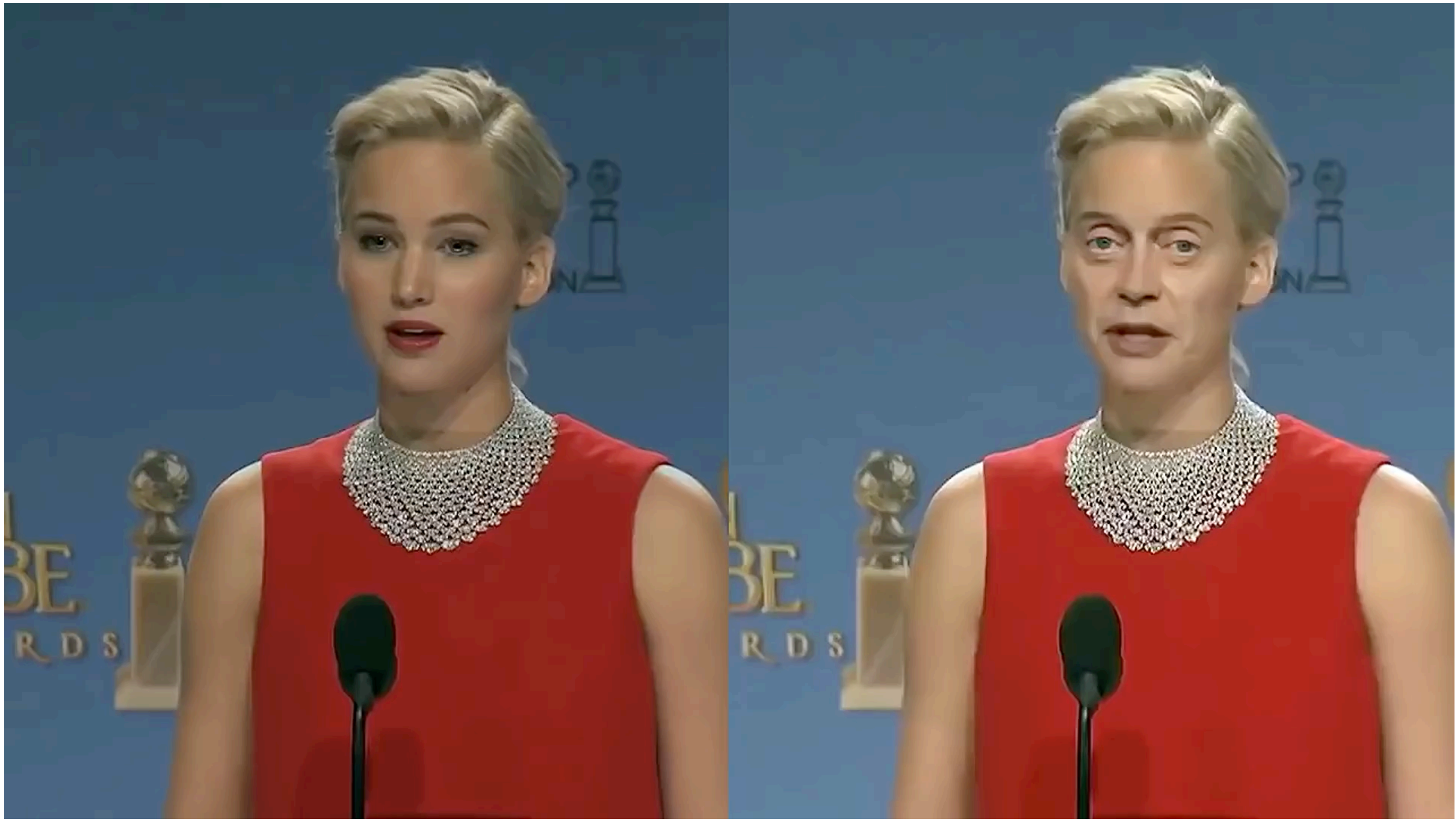
identity



discriminator







TikTok

@deeptomcruise



DeeptomCruise

# This Horrifying App Undresses a Photo of Any Woman With a Single Click

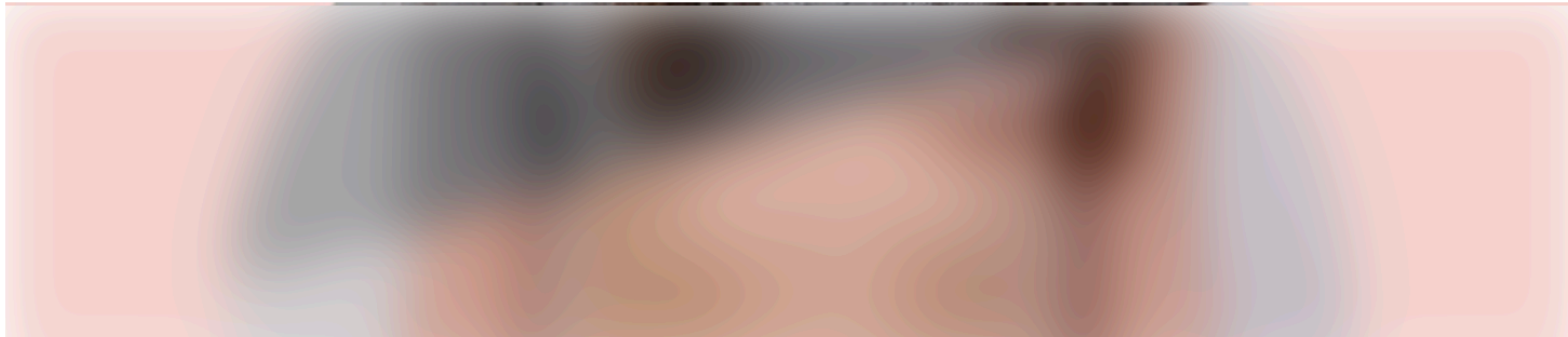
The \$50 DeepNude app dispenses with the idea that deepfakes were about anything besides claiming ownership over women's bodies.



By [Samantha Cole](#)

---

June 26, 2019, 2:48pm  [Share](#)  [Tweet](#)  [Snap](#)



technologists build because they can,

not necessarily because they should

“move fast and break things”

was, is, and always will be a dumb motto

# **(Un)Ethical Computing**