



## **MODULE 12: INTRODUCTION TO AI BIAS & FAIRNESS**

**BA713 - Machine Learning & AI**



## **AI BIAS WITH EXAMPLES**



## **SOURCES & TYPES OF AI BIAS**



## **MEASURING & MITIGATING AI BIAS**

# **CONTENTS**

**Discrimination** of certain individuals, groups, religion, age, gender, skin color, race, culture, economic condition, marital status

**Disparate Impact** → unintentional unfairness that occurs when a decision has widely different outcomes for different groups

AI used to determine important real-world outcomes such as loan approval, pay rates

It is important on AI community to minimize unintentional discrimination

## IMPORTANCE OF REDUCING BIAS IN AI ALGORITHMS

# AI BIAS IS VERY COMMON

## AI expert calls for end to UK use of 'racially biased' algorithms

### Gender bias in AI: building fairer algorithms

### Millions of black people affected by racial bias in health-care algorithms

Study reveals rampant racism in decision-making software used by US hospitals – and highlights ways to correct it.

Overcoming Racial Bias In AI Systems And Startlingly Even In AI Self-Driving Cars

### The Best Algorithms Struggle to Recognize Black Faces Equally

US government tests find even top-performing facial recognition systems misidentify blacks at rates five to 10 times higher than they do whites.

### Bias in AI: A problem recognized but still unresolved

Amazon, Apple, Google, IBM, and Microsoft worse at transcribing black people's voices than white people's with AI voice recognition, study finds

Racial bias in a medical algorithm favors white patients over sicker black patients

## *The Week in Tech: Algorithmic Bias Is Bad. Uncovering It Is Good.*

### AI Bias Could Put Women's Lives At Risk – A Challenge For Regulators

Artificial Intelligence has a gender bias problem – just ask Siri

# WHAT IS AI BIAS?

- An anomaly in the output of an AI algorithm
- Exist in many shapes and forms
- prejudiced assumptions
- Introduced at any stage of the AI pipeline
- **Bias** is present inherently in the world around us, in our society
- We cannot directly solve the bias in the world
- We can take precautions to remove bias from different datasets used to train the AI algorithms
- Since AI has the potential to help humans to make fair decisions
- We need to work towards the **fairness** of the AI algorithms itself
- **Fairness** depends upon the situation in addition to the representation of your values, ethics and legal regulations



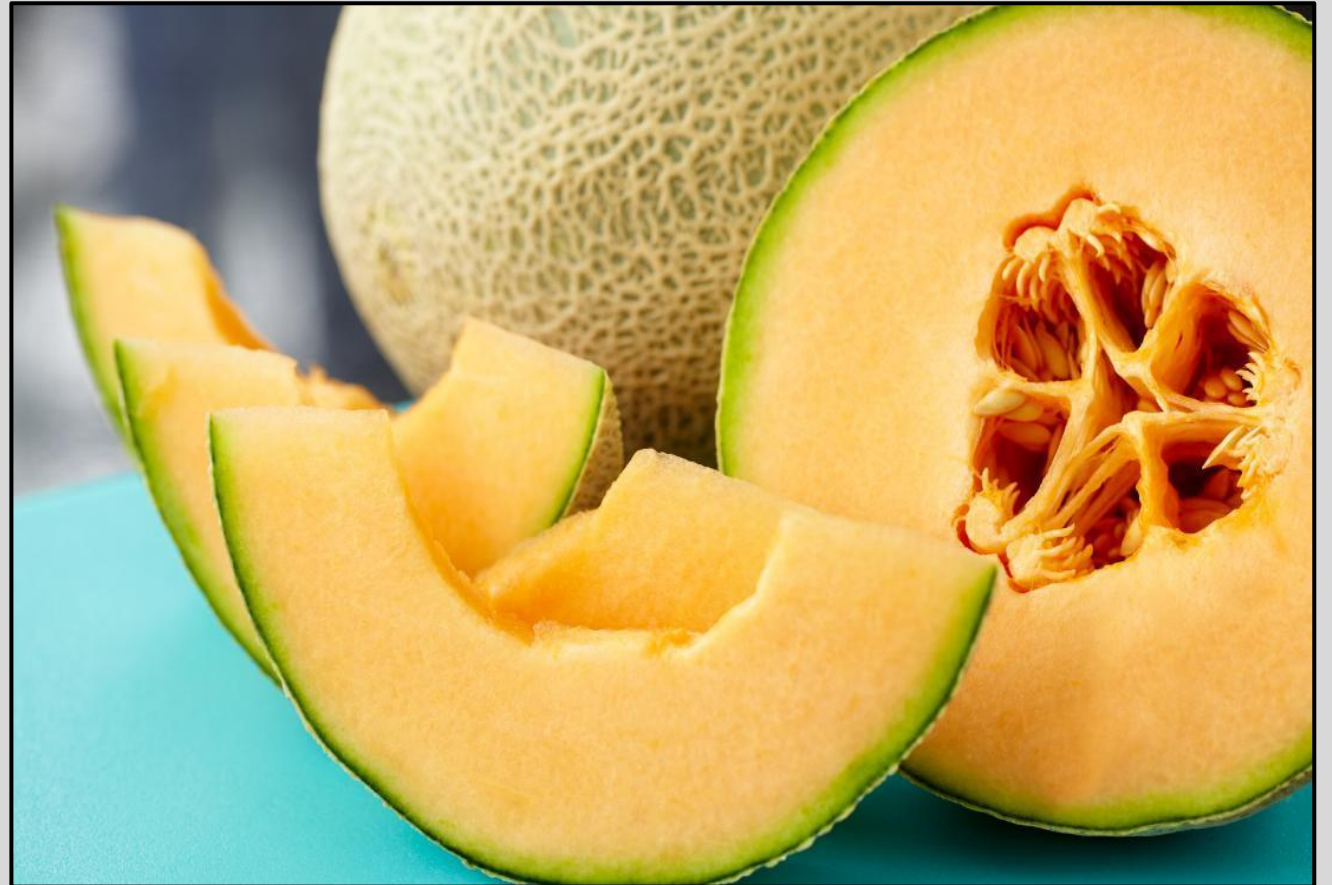
# EXPLAINING BIAS

What do you see in this image?



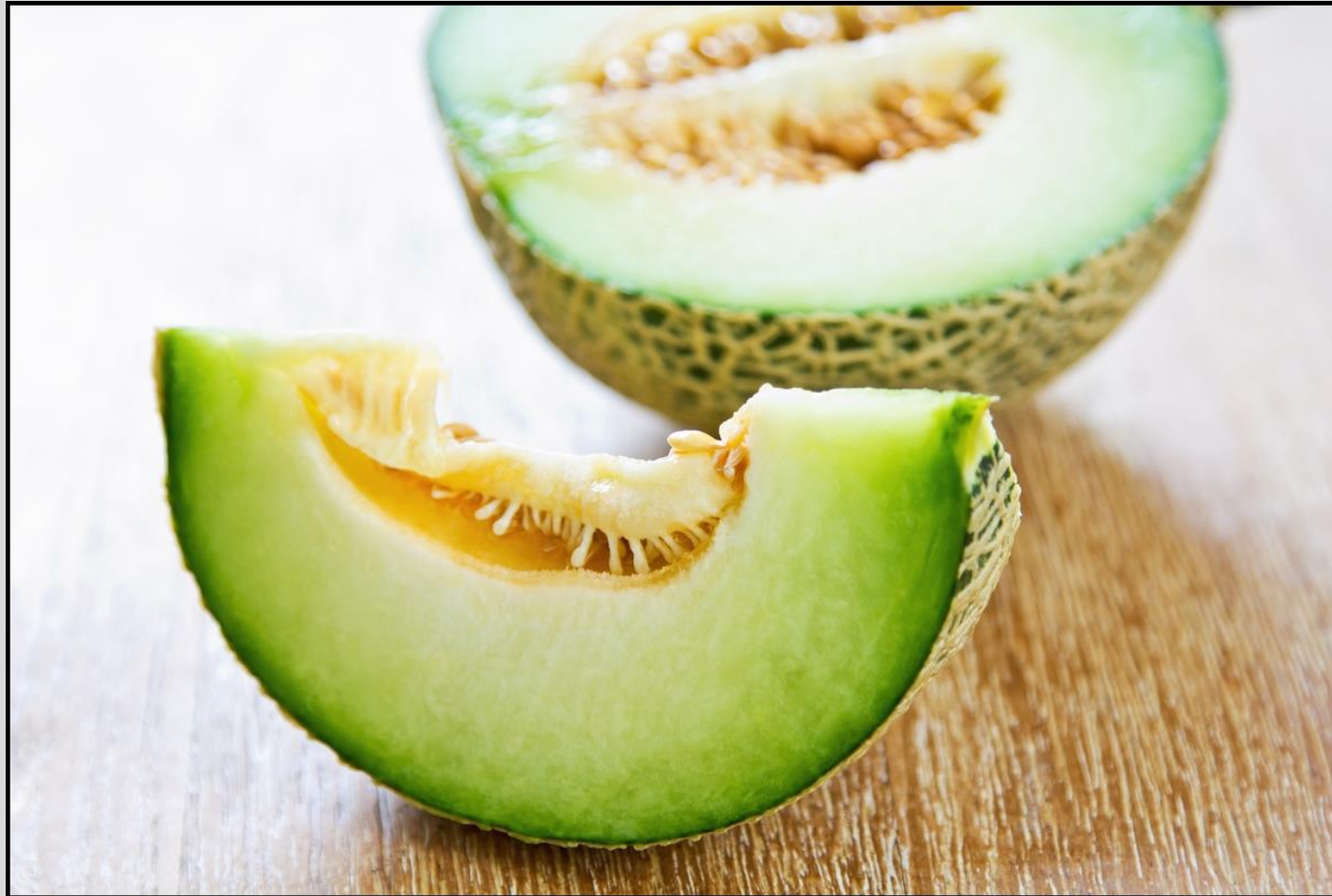
# EXPLAINING BIAS

- cantaloupe (melon)
- cantaloupe slices
- cantaloupe with seeds
- Juicy cantaloupe
- cantaloupe slices next to each other
- But what about **orange cantaloupe** ?



# EXPLAINING BIAS

What do you see in this image?





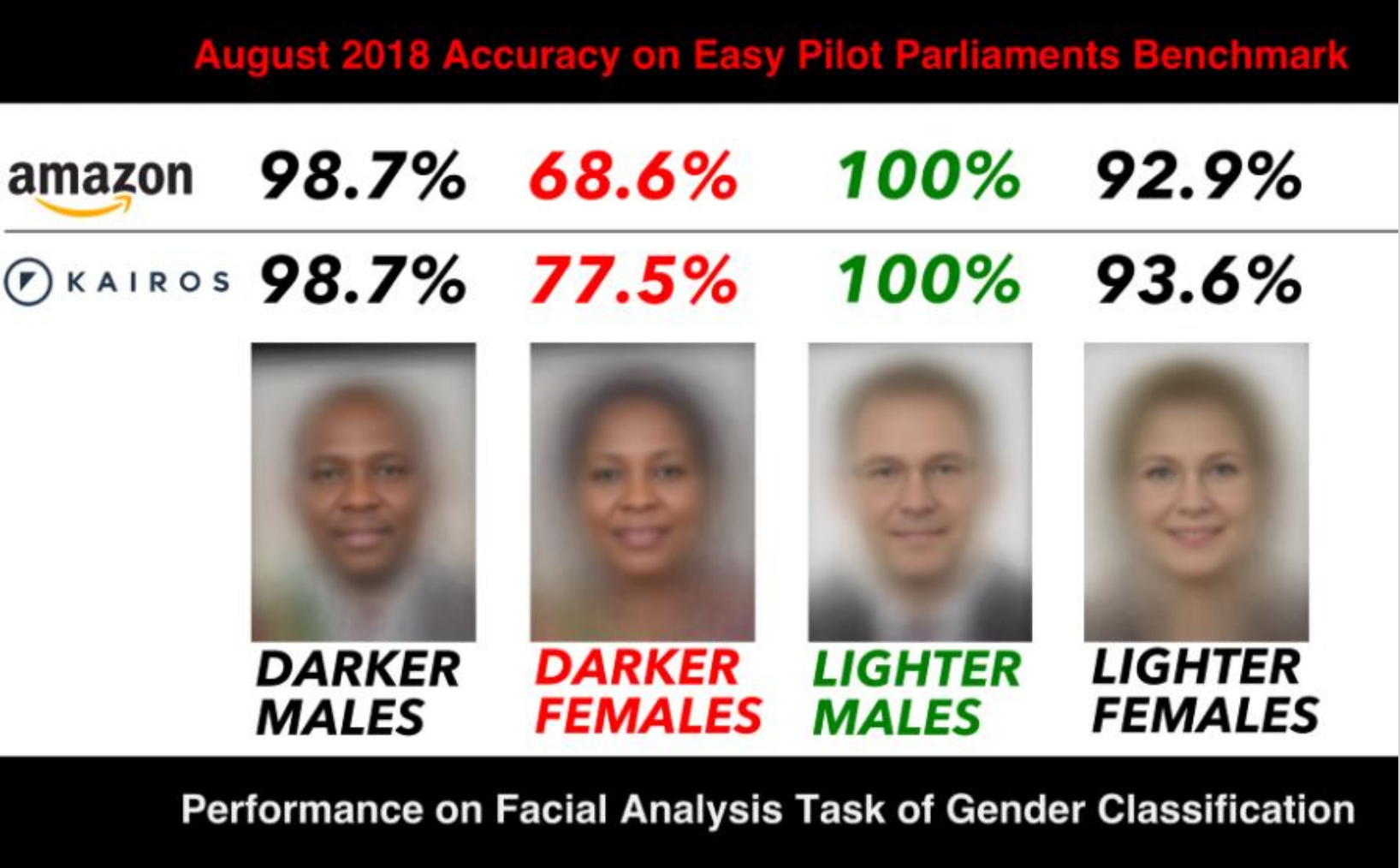
# EXPLAINING BIAS

What do you see in this image?

- **Green** cantaloupe (melon)
- **Green** cantaloupe slices
- **Green** cantaloupe with seeds
- **Green** Juicy cantaloupe
- **Green** cantaloupe slices next to each other

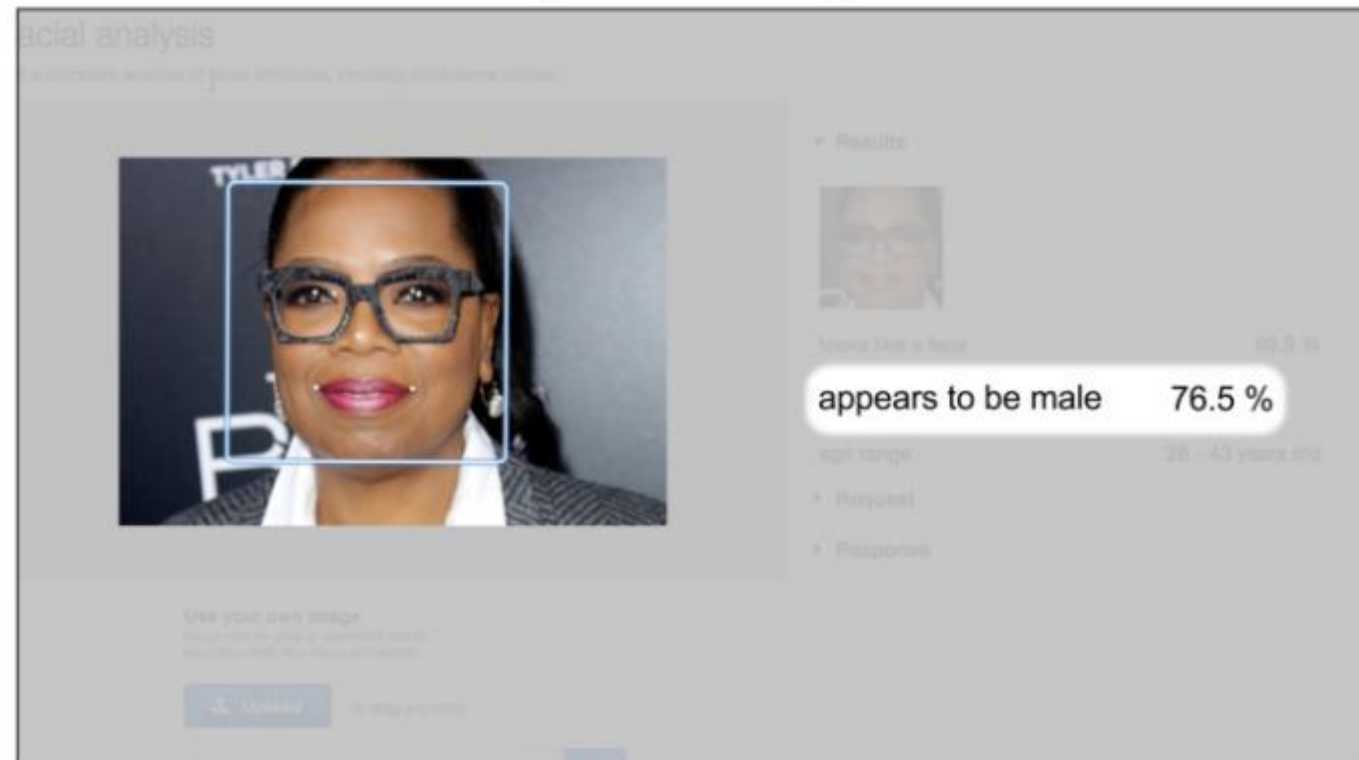


# BIAS IN FACIAL DETECTION



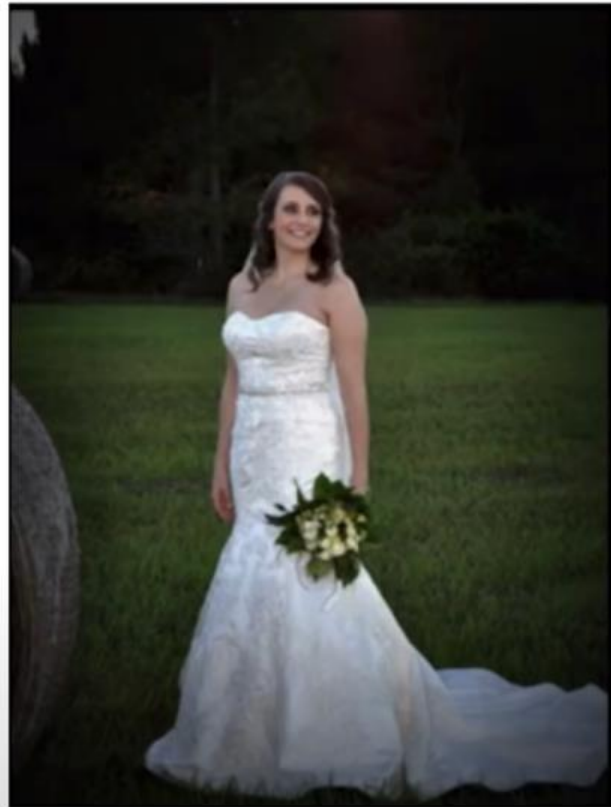
# BIAS IN FACIAL DETECTION

Oprah Winfrey

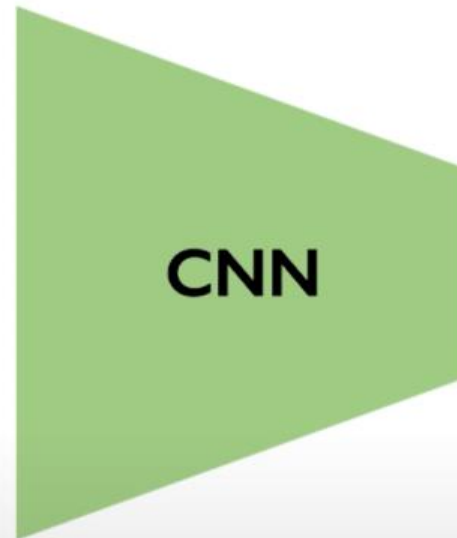


amazon

# BIAS IN IMAGE CLASSIFICATION



Ground Truth: Bride



**CNN**

CNN for image  
classification.



## Predicted Classes

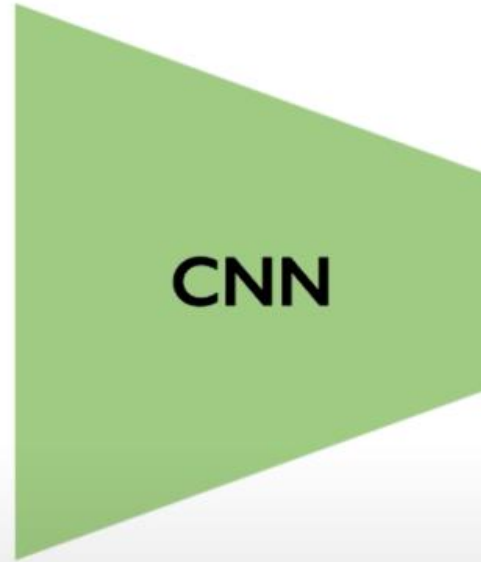
Bride ✓  
Dress  
Ceremony  
Woman  
Wedding



# BIAS IN IMAGE CLASSIFICATION



Ground Truth: Bride



**CNN**

CNN for image classification.



## Predicted Classes

Clothing  
Event  
Costume  
Red  
Performance art



# REAL-LIFE EXAMPLES OF AI BIAS

## **COMPAS:**

software used by the US courts to judge the probability of a defendant → Biased towards black defendants

## **Amazon's AI Recruiting System (2014):**

discriminatory against women, use data from past 10 years, most selected applicants were men due to the male dominance in the tech industry, scrapped in 2018

## **US Healthcare:**

Since it is perceived, that Black people are less able to pay the costs, the AI ranked their health risk lower than white people for the same disease

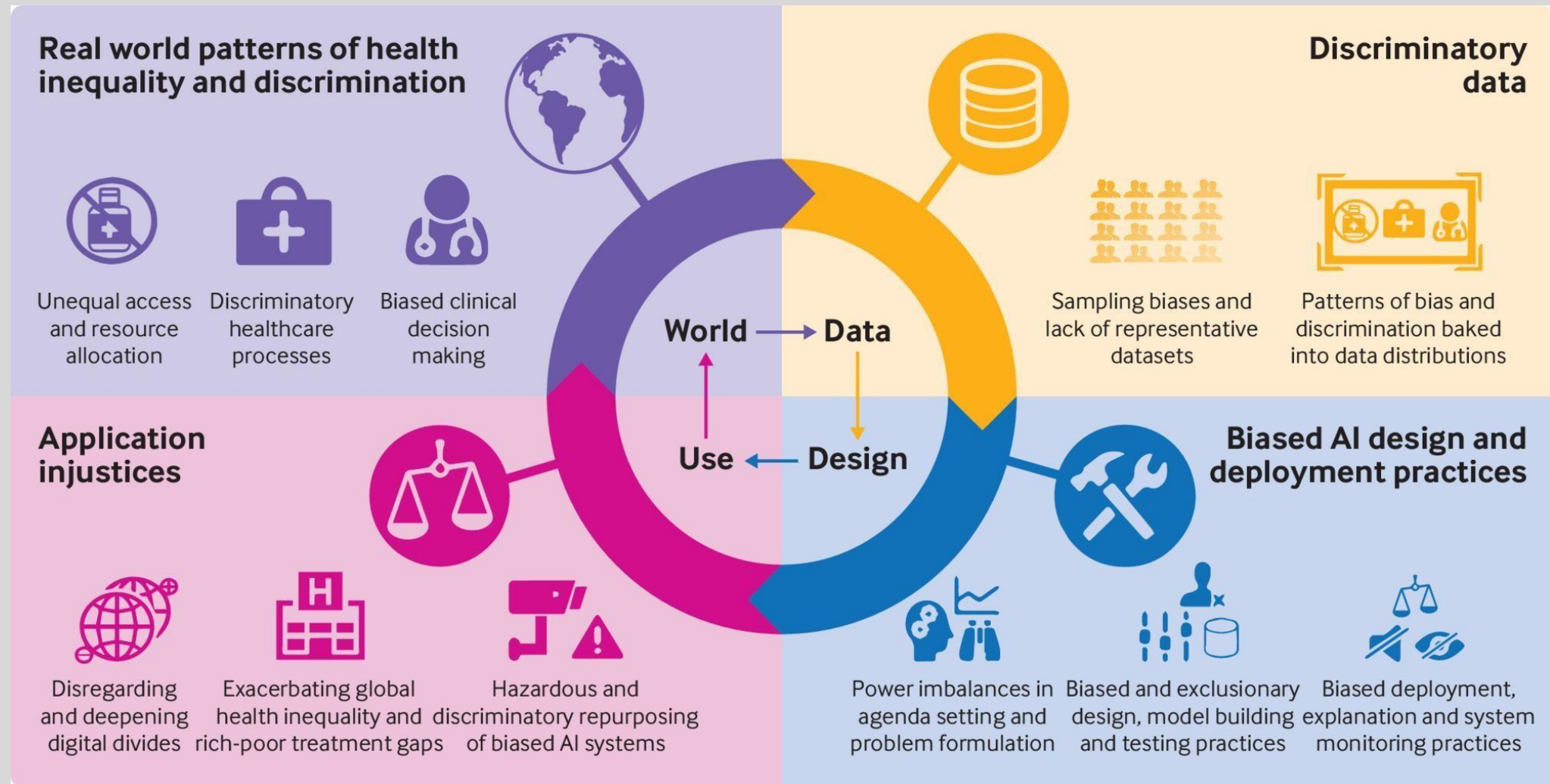
## **Twitter Image Cropping (2020):**

image cropping algorithm favored white faces over black faces

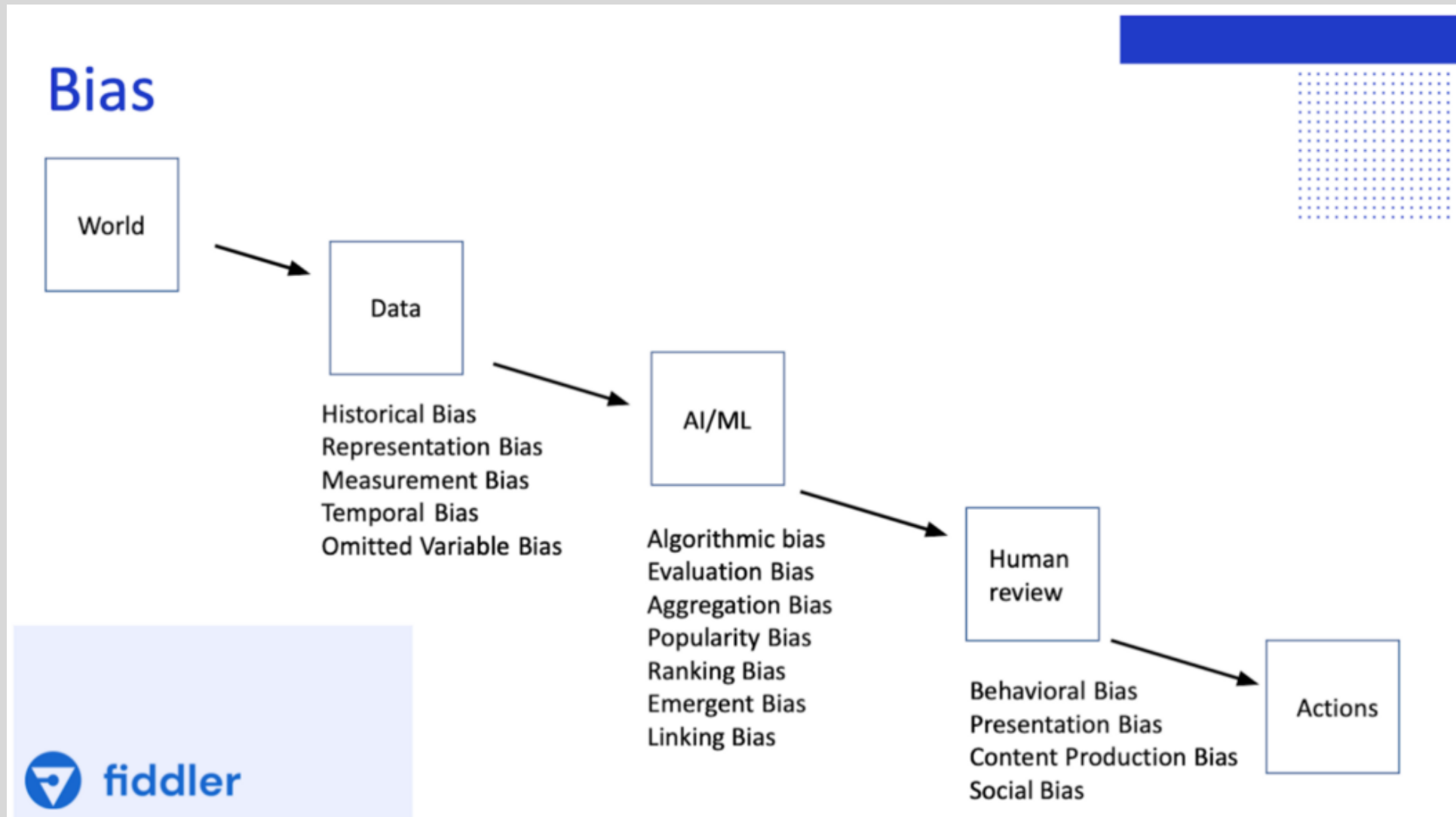
## **Facebook's Advertisement Algorithm (2019):**

advertisers target people based on their race, gender, and religion  
For example, jobs like nursing and secretary being targeted to women while jobs like janitor and taxi drivers targeted to men of color

# SOURCES OF AI BIASES



# TYPES OF AI BIAS





# BIAS IN THE AI TRAINING DATASETS

- Data imbalance, class labels, features, input structure

## 1) Historical Bias:

- already exists in the world and is present in our training data
- **For example**, 2016 paper “Man is to Computer Programmer as Woman is to Homemaker”

## 2) Representation bias:

- the way we define and sample a population to create a dataset
- **For example**, the data used to train Amazon's facial recognition was mostly based on white faces

## 3) Measurement bias:

- When collecting features or labels to use in predictive models
- **For example**, predictive policing, likelihood that a black defendants will commit another crime leading to harsher sentences for them as compared to white defendants

# BIAS IN AI MODELS

- uncertainty, interpretability, and performance metrics, training feedback loops perpetuate biases
- lack of systematic analysis with respect to data subgroups

## 1) Evaluation Bias:

- Occurs during model iteration and evaluation
- Benchmark models do not represent the general population, or are not appropriate for the way the model will be used
- **For example**, evaluation metrics such as balanced accuracy and FPR must be consistent for different social groups (gender, ethnicity or age)

## 2) Aggregation bias:

- Occurs during model construction
- distinct populations are inappropriately combined
- For the population is heterogenous, a single model will not suit the needs of all groups
- **For example**, in medical applications such as diagnosing and monitoring diabetes, AI models use Hemoglobin levels as a predictor. But research studies show that these levels may differ across different ethnic groups

# BIAS IN HUMAN REVIEW

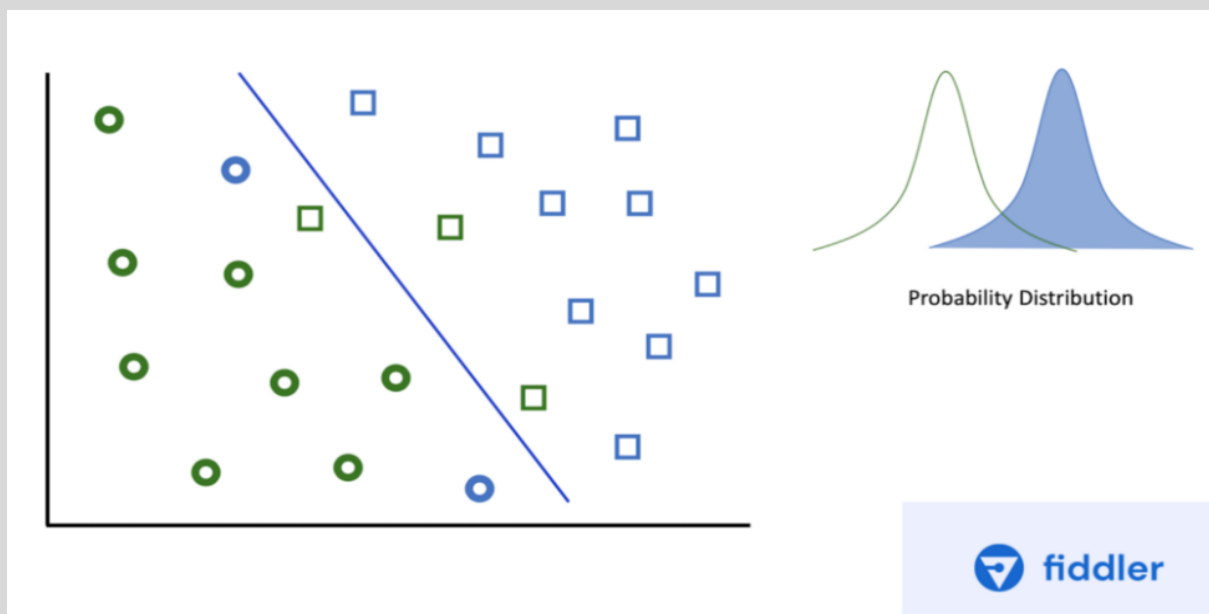
- human errors and biases distort the meaning of AI predictions

## Social Bias:

- A human reviewer may introduce their own biases when accepting or disregarding the AI model's predictions
- **For example**, a human reviewer might override a correct prediction done by the AI model based on their social bias. **"I know since this defendant is black, he will commit another crime"**

# WHAT IS FAIRNESS?

- **Fairness:** absence of prejudice or preference for an individual or group based on their characteristics
- **For example,** in the figure, data included two different underlying groups → green and blue
- represent different ethnicities, genders, or even geographical or temporal differences like morning and evening users
- In this case, using a single threshold for these two groups of diabetes patients would lead to poor health outcomes
- predictions are calibrated for each group properly





### Fairness of error:

- similar performance (**accuracy**) across different groups, no group has more error than the other

### Fairness of representation:

- model is equitable in its assignment of favorable outcomes to each group
- **Proportional parity**: metric → requires that each group receives the favorable outcome the same proportion of the time
- **For example**, female and male candidates in automated resume review are each accepted 50% of the time

## MEASURING AI BIAS

# MEASURING FAIRNESS

## Predictive parity:

- Measures the raw accuracy score between groups
- Number of times the outcome is correctly identified

## False positive and negative rate parity:

- measures the model's false positive rate and false negative rate for each class



**AWARENESS &  
UNDERSTANDING OF  
DATA & AI  
ALGORITHM**



**IMPROVING DATA COLLECTION,  
RESAMPLING DATA USING  
GENERATIVE MODELS, REDUCING  
CLASS IMBALANCE PROBLEM**



**IMPROVING HUMAN  
CENTRIC DESIGN  
APPROACH**



**ESTABLISH A DEBIASING  
STRATEGY SUCH AS,  
EVALUATION METRICS,  
SUBGROUPS AND  
COMBINATIONS**



**MAINTAINING A  
DIVERSE AI TEAM,  
ENGAGE FACT-BASED  
CONVERSATIONS  
ABOUT POTENTIAL AI  
BIASES**



**MONITORING AND  
UPDATING YOUR  
MODEL DURING  
TRAINING & AFTER  
DEPLOYMENT**



**INVEST MORE IN BIAS  
RESEARCH, MULTI-  
DISCIPLINARY  
APPROACH**

**Minimizing AI bias is important for AI systems to flourish and increase people's trust in such systems**

# MITIGATING AI BIAS

# AUTOMATIC RESAMPLING USING GENERATIVE MODELS

- Use latent distribution to generate new data
- Reduce AI biases



Homogeneous skin color, pose

VS

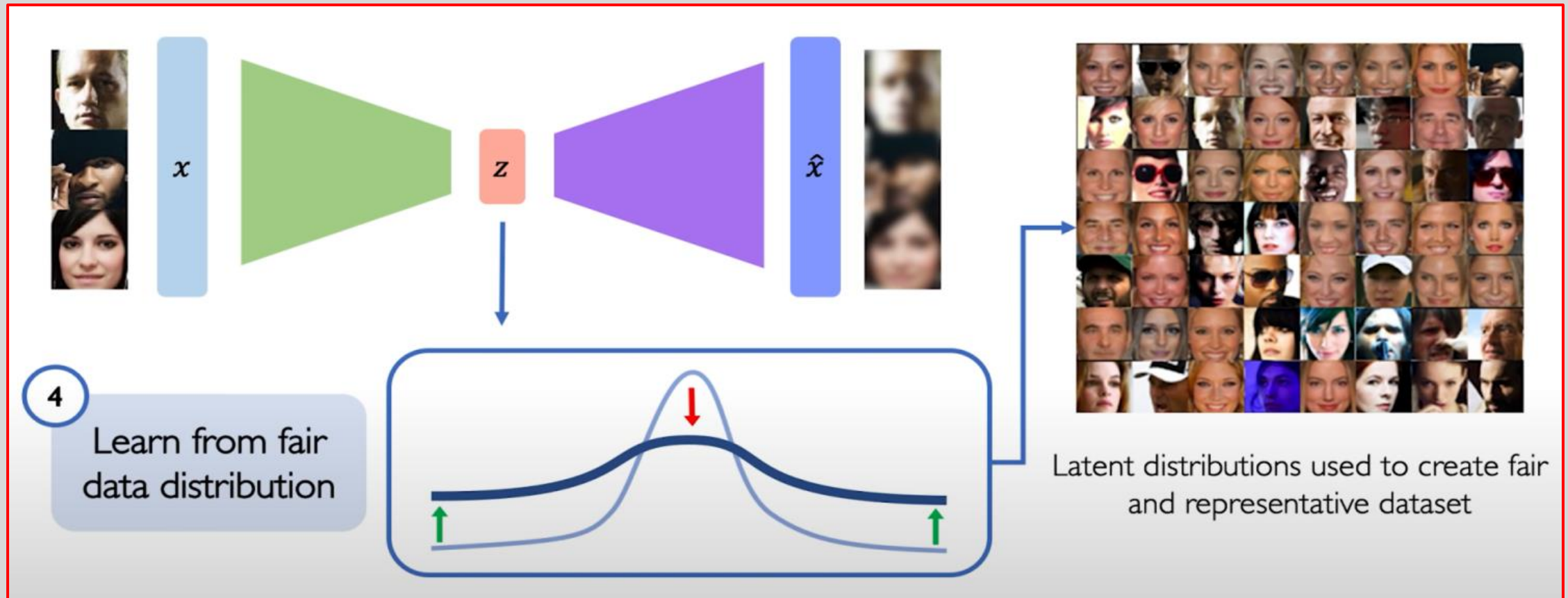


Diverse skin color, pose, illumination



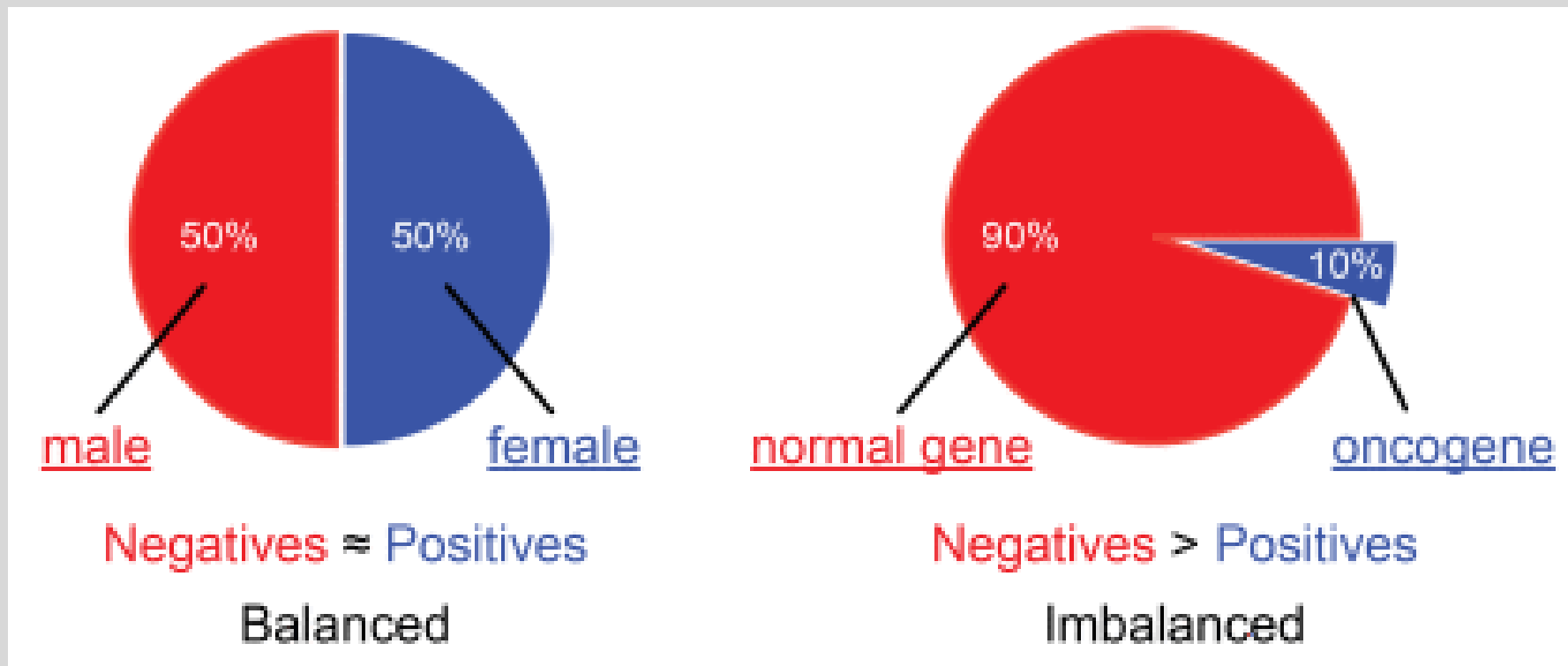
# MITIGATING DATASET BIAS USING VAE

- Generate a fairer and more representative dataset for training



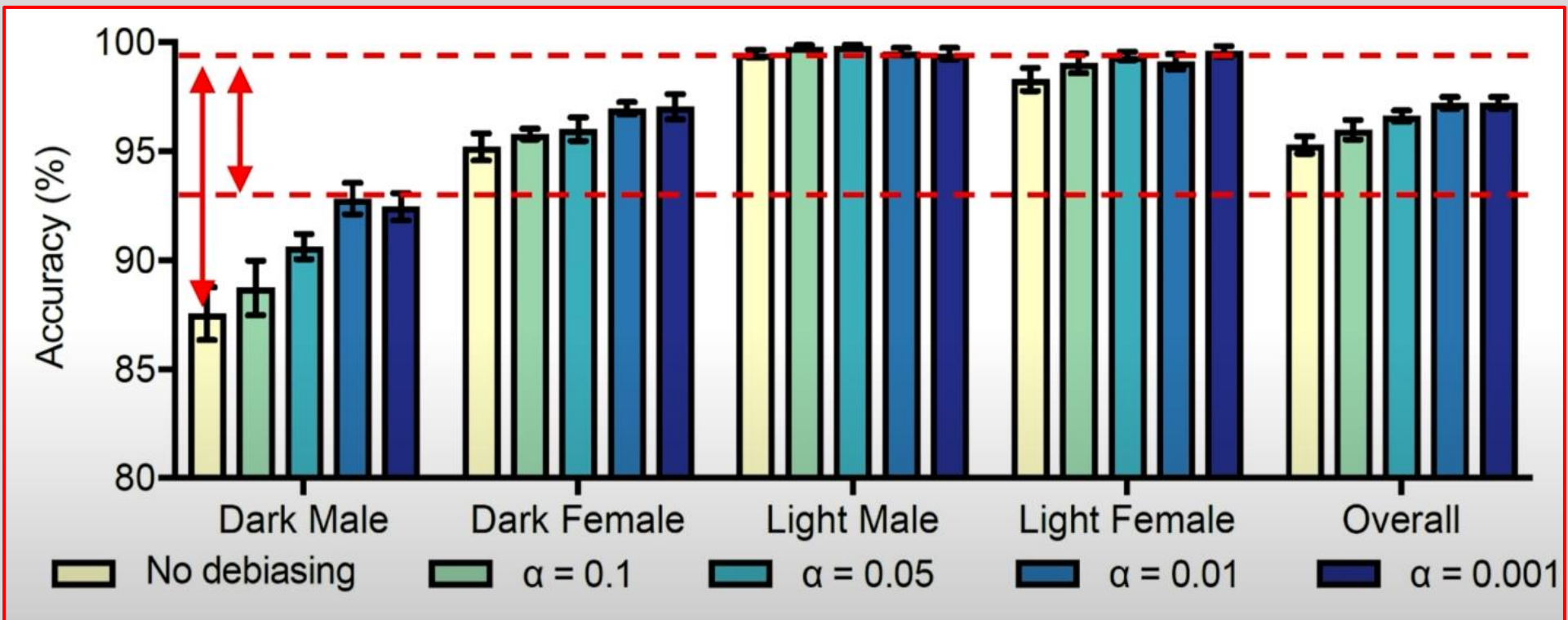
# BALANCED DATASET FOR EVALUATION

- **Balanced Dataset:** equal or approximate number of samples from class A and B
- **Oversampling** and **Undersampling** techniques may be employed



# EVALUATION TECHNIQUE TO REDUCE CATEGORICAL BIAS

- Decrease categorical bias
- Evaluate across subgroups and combinations of subgroups
- Disaggregated and intersectional evaluation



# TOOLS TO REDUCE AI BIAS & DEVELOP A FAIR AI

## AI Fairness 360:

- IBM
- Open source
- detect and mitigate biases in unsupervised learning
- test biases in models and datasets with a comprehensive set of metrics
- binary classification problems only

## IBM Watson OpenScale:

- IBM
- bias checking and mitigation in real time
- when AI is making its decisions.

## Google's What-If Tool:

- can test performance in hypothetical situations
- analyze the importance of different data features
- visualize model behavior across multiple models and subsets of input data
- for different ML fairness metrics

## FATE:

- Fairness, Accountability, Transparency, and Ethics in AI (FATE)
- Microsoft
- assess visualization dashboards and bias mitigation algorithms
- compare trade-offs between fairness and performance of the system



**THANKS!**

**Do you have any  
questions?**