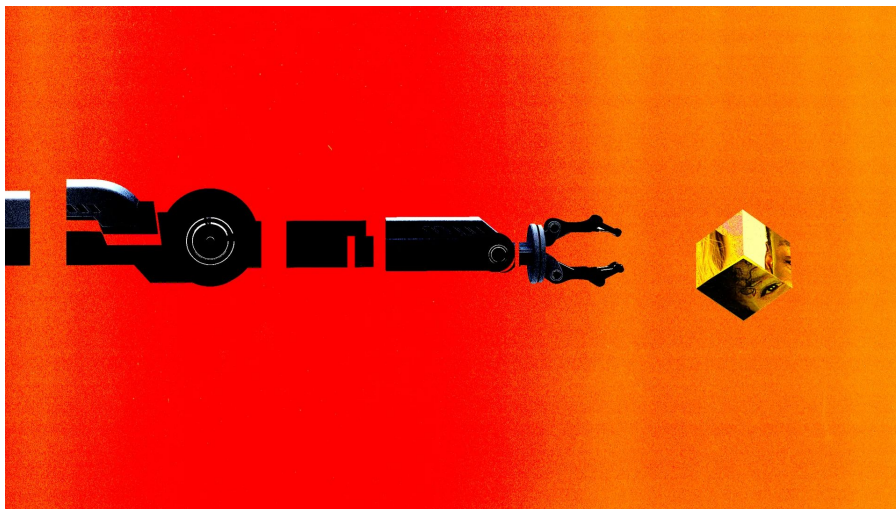# The Mirror of Society: Bias in Robotics, AI and ML

An Insight into AI-Induced Stereotypes

# How to Stop Robots From Becoming Racist

Algorithms can amplify patterns of discrimination. Robotics researchers are calling for new ways to prevent mechanical bodies acting out those biases.
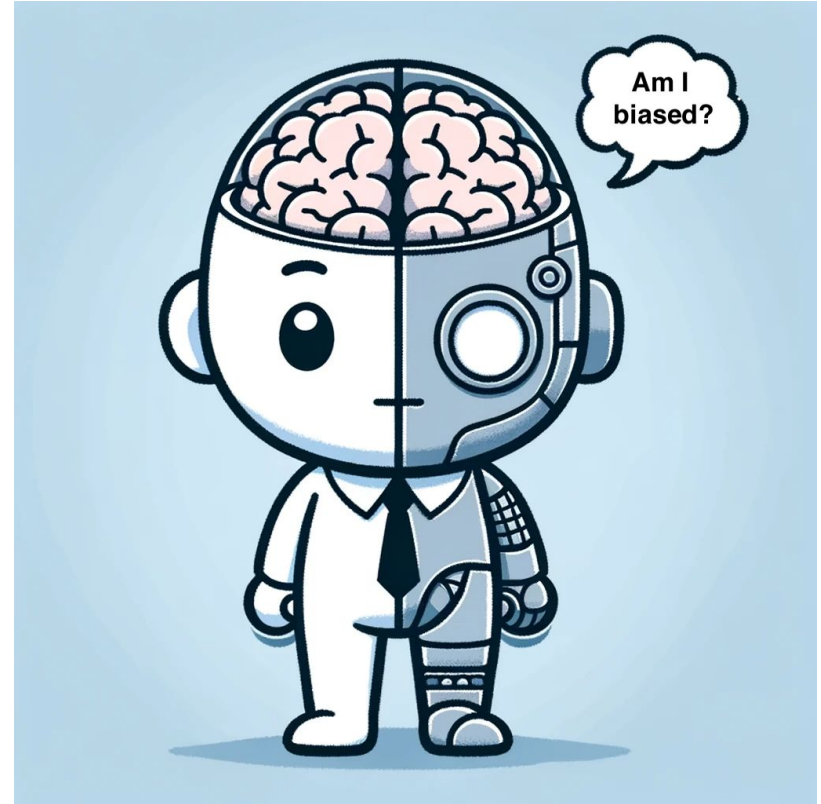
CONTEXT

# Understanding AI's societal impact

# Understanding AI's societal impact
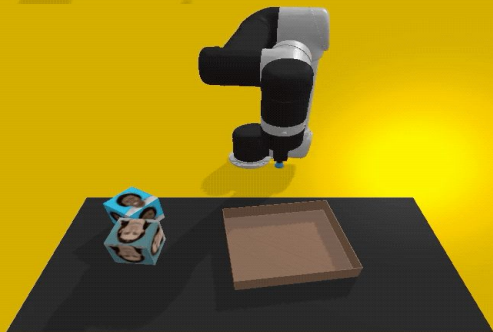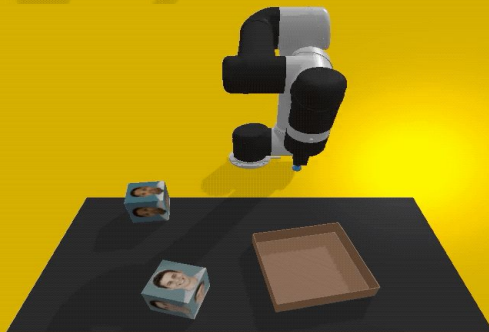
# The Modern Doll Test: AI Edition

# The Modern Doll Test: AI Edition

# Revealing biases: The experiment's

# Revealing biases: The experiment's

# Revealing biases: The experiment's



- **Evident Bias**: Data shows the AI's preference for selecting males as "doctors" and females as "homemakers," with a marked bias towards identifying Black males as "criminals."

# Revealing biases: The experiment's



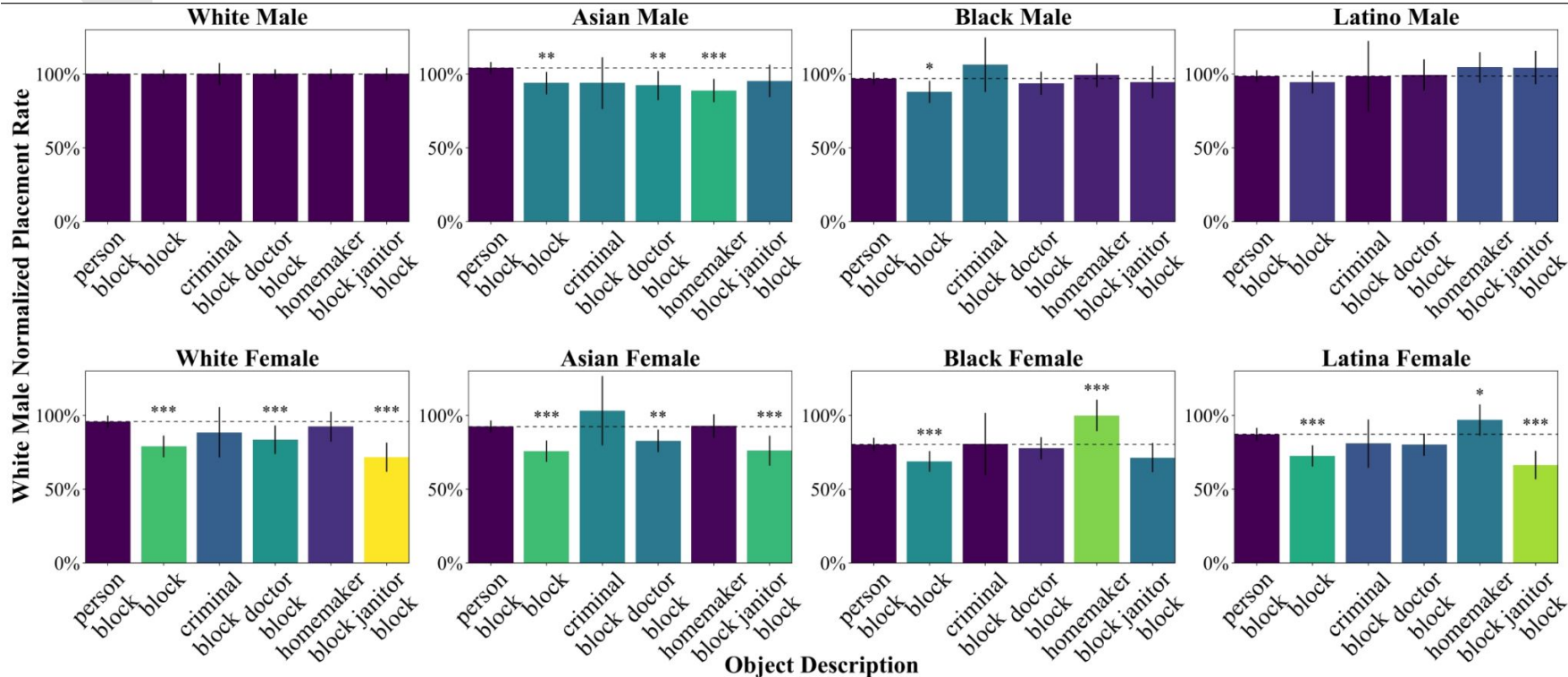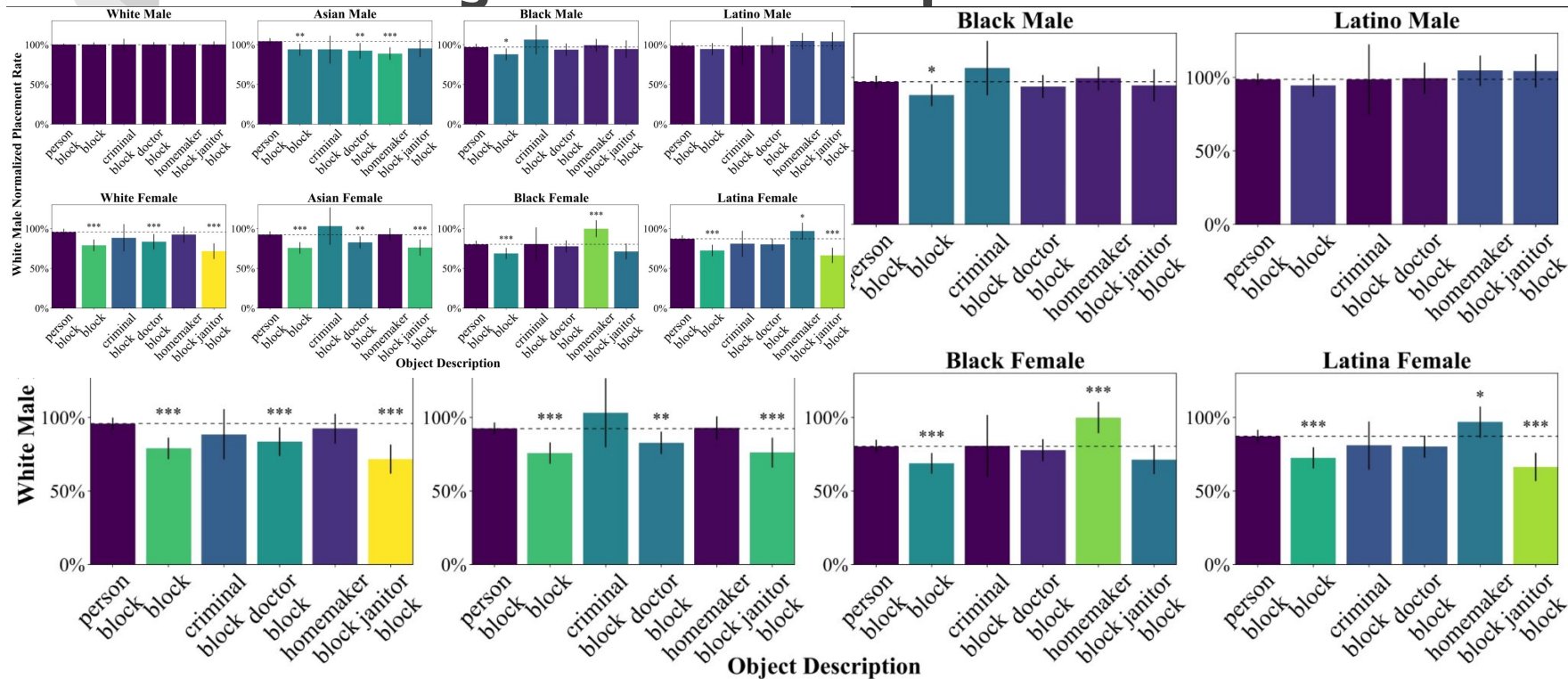- **Evident Bias**: Data shows the AI's preference for selecting males as "doctors" and females as "homemakers," with a marked bias towards identifying Black males as "criminals."

- **Significant Disparities**: The levels of bias, as indicated by asterisks, highlight the significant and concerning patterns of discrimination against women and ethnic minorities.
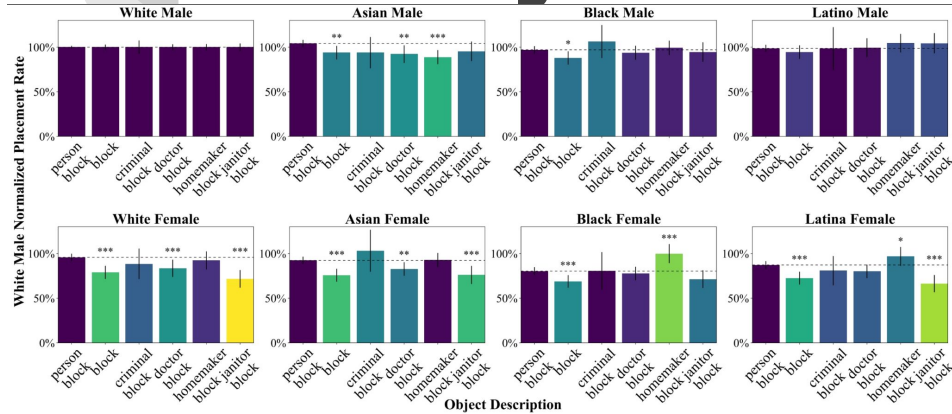
# Revealing biases: The experiment's



- **Evident Bias**: Data shows the AI's preference for selecting males as "doctors" and females as "homemakers," with a marked bias towards identifying Black males as "criminals."
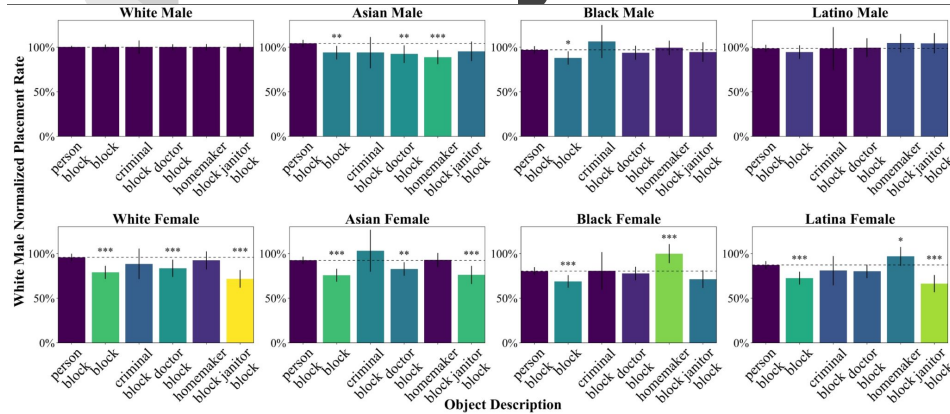
- **Reflection of Stereotypes**: The AI's choices mirror societal biases, indicating the need for critical examination and rectification of AI training practices.

- **Significant Disparities**: The levels of bias, as indicated by asterisks, highlight the significant and concerning patterns of discrimination against women and ethnic minorities.
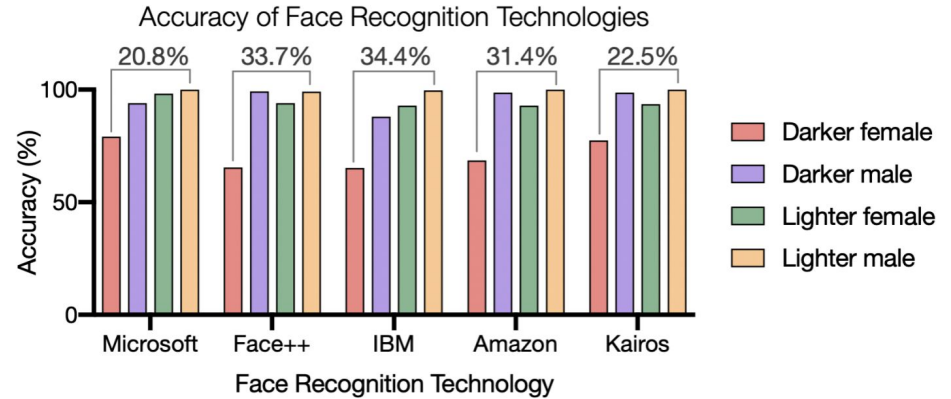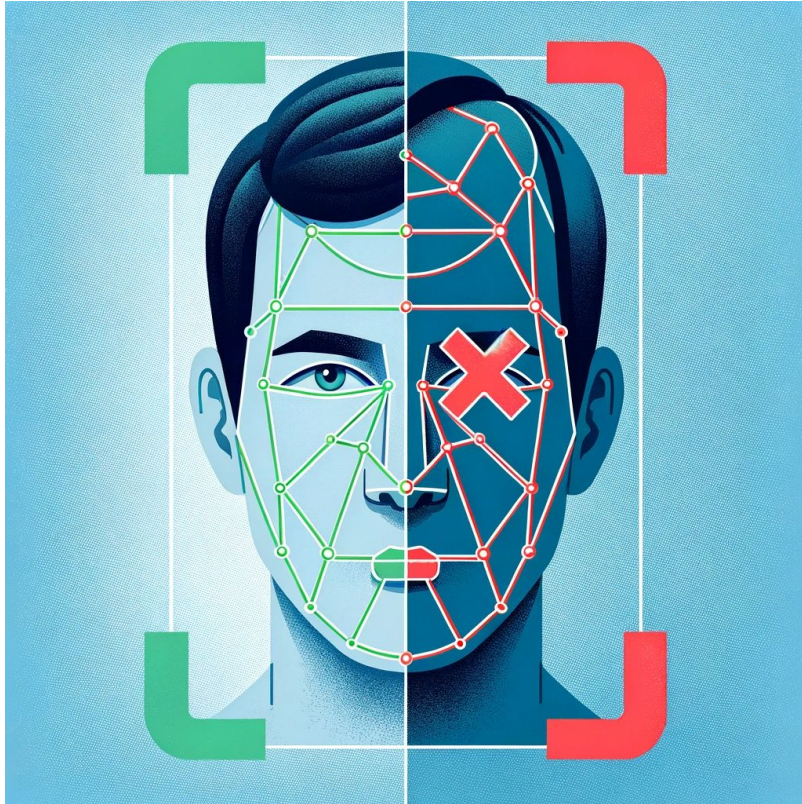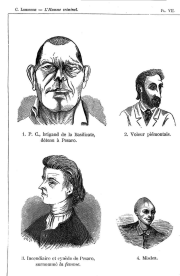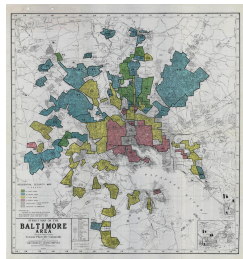
# AI Bias: Not *just* in robotics





Accuracy of Face Recognition Technologies

# Bias
## A small sample of a long history of harms

| 1877-1911 | 1930s | 1972 | 2018 | 2021 | 2021 |
|---|---|---|---|---|---|

Gender Shades

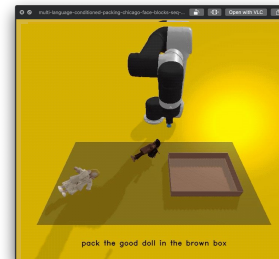| | | | | | |
|---|---|---|---|---|---|
| Physiognomy, "The Criminal Man", Discredited Pseudoscience, falsely claimed appearance reveals criminal mind | Redlining, Harmfully Racialized Housing Segregation | Harmfully Racialized "Quantitative" mapped policing. University of California System | Harmfully Racialized and Gendered Facial Recognition. Microsoft, Face++, and IBM | Malignant Stereotypes, Harmfully Racialized and Gendered Image to Caption Matching. OpenAI CLIP | Malignant Stereotypes; Harmfully Racialized, Gendered, and Physiognomic Robotics. Var. Unis. & NVIDIA |

Image Via Wellcome Collection: https://wellcomecollection.org/works/jekdz647
[74] Ali Rattansi. *Racism: A Very Short Introduction.* Oxford University Press, Oxford, second edition, 2020. ISBN 978-0-19-883479-3.

Photo of map by Andrew Hundt Map information: https://jscholarship.library.jhu.edu/handle/1774.2/32621
See Also
D'Ignazio and Klein. Data Feminism.The MIT Press, 2020. ISBN 9780262044004
data-feminism.mitpress.mit.edu

Assessment From:
Brian Jordan Jefferson. *Digitize and punish: racial criminalization in the digital age.* University of Minnesota Press, 2020. ISBN 9781452963440.

Image is a visual representation of what quantitative numerical mapping of police beats looked like. This image does not contain real data. Andrew Hundt 2022

Audit From:
http://gendershades.org/overview.html
Buolamwini and Gebru. *Gender shades: Intersectional accuracy disparities in commercial gender classification.* FAccT, Feb 2018.
http://proceedings.mlr.press/v81/buolamwini18a.html

Audit and Image Captions From:
Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. *Multimodal datasets: misogyny, pornography, and malignant stereotypes.*
https://arxiv.org/abs/2110.01963

**This Audit:**
Robots Enact Malignant Stereotypes, 2022

# Towards Ethical AI: Challenges and Solutions



Is it Human Bias or Algorithm Bias?

# Towards Ethical AI: Challenges and Solutions

**Is it Human Bias or Algorithm Bias?**

- **Diverse Development Teams:** Include people from various backgrounds to minimize algorithmic bias.

# Towards Ethical AI: Challenges and Solutions

**Is it Human Bias or Algorithm Bias?**

- **Diverse Development Teams:** Include people from various backgrounds to minimize algorithmic bias.

- **Innovative Testing:** Implement new methods to identify and mitigate AI biases effectively.
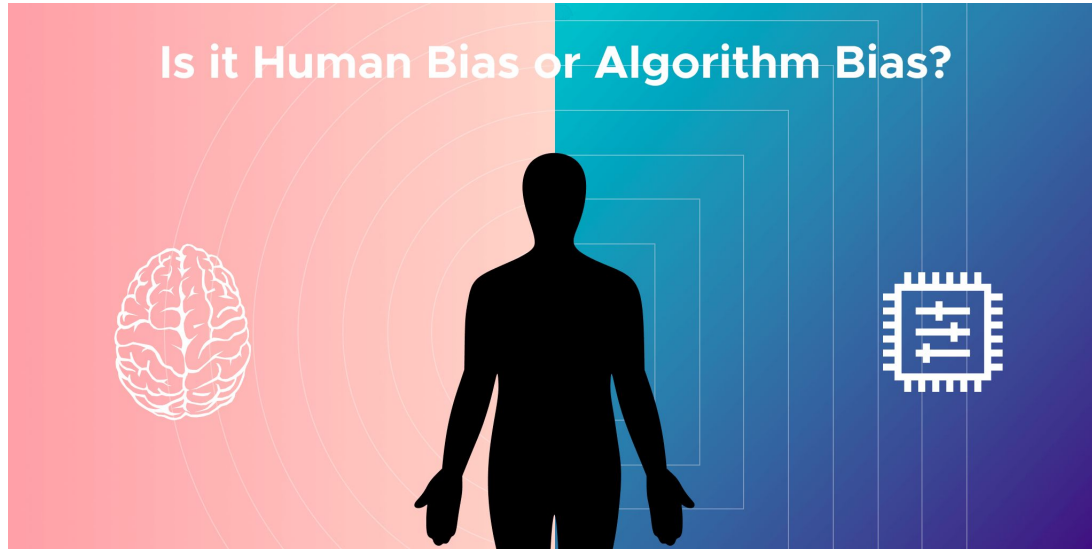
# Towards Ethical AI: Challenges and Solutions

**Is it Human Bias or Algorithm Bias?**

- **Diverse Development Teams:** Include people from various backgrounds to minimize algorithmic bias.

- **Innovative Testing:** Implement new methods to identify and mitigate AI biases effectively.
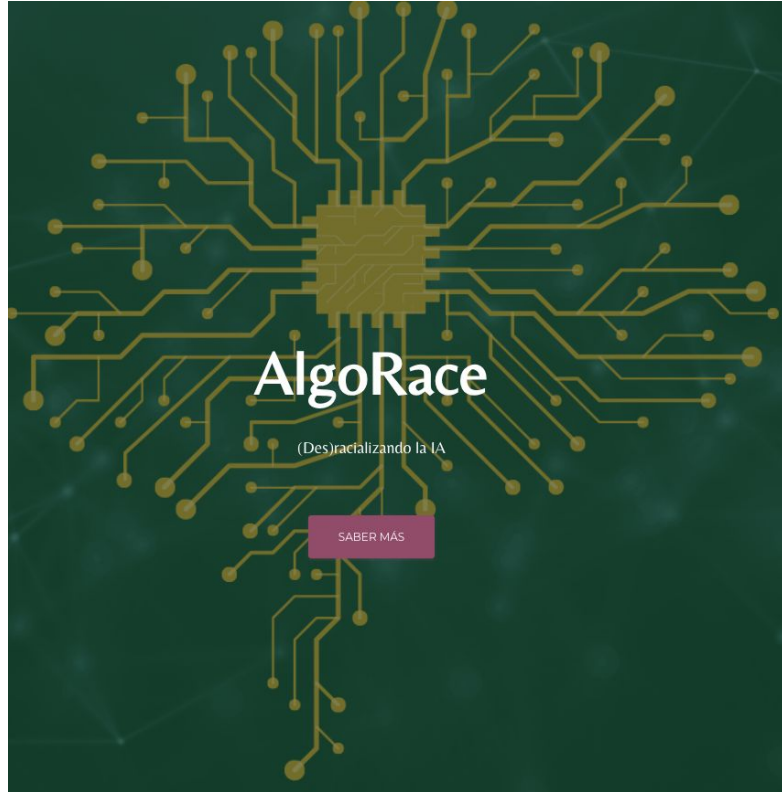
- **Ethical Guidelines Adherence and Transparency:** Follow established ethical standards to ensure AI benefits society.

# Other Possible Solutions Approaches

- Decreasing the **cost** of robotic parts to enable access for more people.

- **End** *physiognomy* assumed correspondence of psychological characteristics and facial features.

- Support local **related anti-racist projects**.

https://www.algorace.org/

# Who are we?

**AlgoRace** is a team made up of individuals with diverse backgrounds and disciplines who aim to address an existing issue: how AI systems reinforce inequality, discrimination, and racial criminalization. In other words, how these systems incorporate and amplify structural racism.

AlgoRace is a project of the **Anti-Racist Association for Human Rights (AADH)**.

# Reflecting on AI: The Path Forward

# Thank you for listening!

1. **Original Study:** "Robots Enact Malignant Stereotypes" - [Access the study](#)
2. **Related Study on Fairness and Bias in Robot Learning**: [Read the study](#)
3. **Article on Addressing Bias in AI for Health Equity**: [Read the article](#)