

Ethics

Preamble

- This is a huge subject, spanning many disciplines and addressing many real different problems.
- We, ML practitioners, need to have a student mindset, and not assume we have the answers. These are not easy problems.
- I am not an expert. Excellent resources recommended at the end.

Tech Ethics Curricula

Table 2: The number of courses that had content for each listed topic, out of 115 total courses, organized from most popular to least popular topic.

| Topic | Courses |
|---------------------------------------|---------|
| Law & policy | 66 |
| Privacy & surveillance | 61 |
| Philosophy | 61 |
| Inequality, justice & human rights | 59 |
| AI & algorithms | 55 |
| Social & environmental impact | 50 |
| Civic responsibility & misinformation | 32 |
| AI & robots | 27 |
| Business & economics | 27 |
| Professional ethics | 25 |
| Work & labor | 23 |
| Design | 20 |
| Cybersecurity | 19 |
| Research ethics | 16 |
| Medical/health | 12 |

Table 3: The number of courses that had each type of learning outcome, organized from most common to most common outcome.

| Outcome | Courses |
|---------------------------|---------|
| Critique | 71 |
| Spot issues | 36 |
| Make arguments | 26 |
| Improve communication | 26 |
| See multiple perspectives | 23 |
| Create solutions | 21 |
| Consider consequences | 18 |
| Apply rules | 10 |

What Do We Teach When We Teach Tech Ethics? A Syllabi Analysis

Casey Fiesler
casey.fiesler@colorado.edu
University of Colorado Boulder
Boulder, CO

Natalie Garrett
natalie.garrett@colorado.edu
University of Colorado Boulder
Boulder, CO

Nathan Beard
nbeard@umd.edu
University of Maryland
College Park, MD

SIGCSE 2020

- 115 university tech ethics courses analyzed
- Huge variation in materials
- More consensus in outcomes: ability to critique, spot issues, and make and communicate arguments.

One last preamble

There are these two young fish swimming along, and they happen to meet an older fish swimming the other way, who nods at them and says, “Morning, boys. How’s the water?” And the two young fish swim on for a bit, and then eventually one of them looks over at the other and goes, “What the hell is water?”

<https://www.newyorker.com/books/page-turner/this-is-water>

Outline

- What is ethics
- Long-term ethical problems in AI
- Near-term ethical problems in AI
- Best practices
- Where to learn more

What is ethics

Ethics?

Ethics ≠ Feelings

Ethics ≠ Laws

Ethics ≠ Societal Beliefs

<https://www.scu.edu/ethics/ethics-resources/ethical-decision-making/what-is-ethics/>

Ethical Theories

Divine Command

- Moral behaviors are those commanded by the divine
- Criticism: not much philosophy can say

Virtue Ethics

- Moral behaviors uphold the person's virtues
- Criticism: increasing evidence that character traits are illusory

Deontology (Duty)

- Moral behaviors are those that satisfy the categorical imperative (e.g. don't lie, don't kill)
- Criticism: unacceptable inflexibility

Utilitarianism

- Moral behaviors are those that bring the most good to the most people
- Criticism: How to measure utility?

<https://kevinbinz.com/2017/04/13/ethical-theory-intro/>

Is one a clear winner?

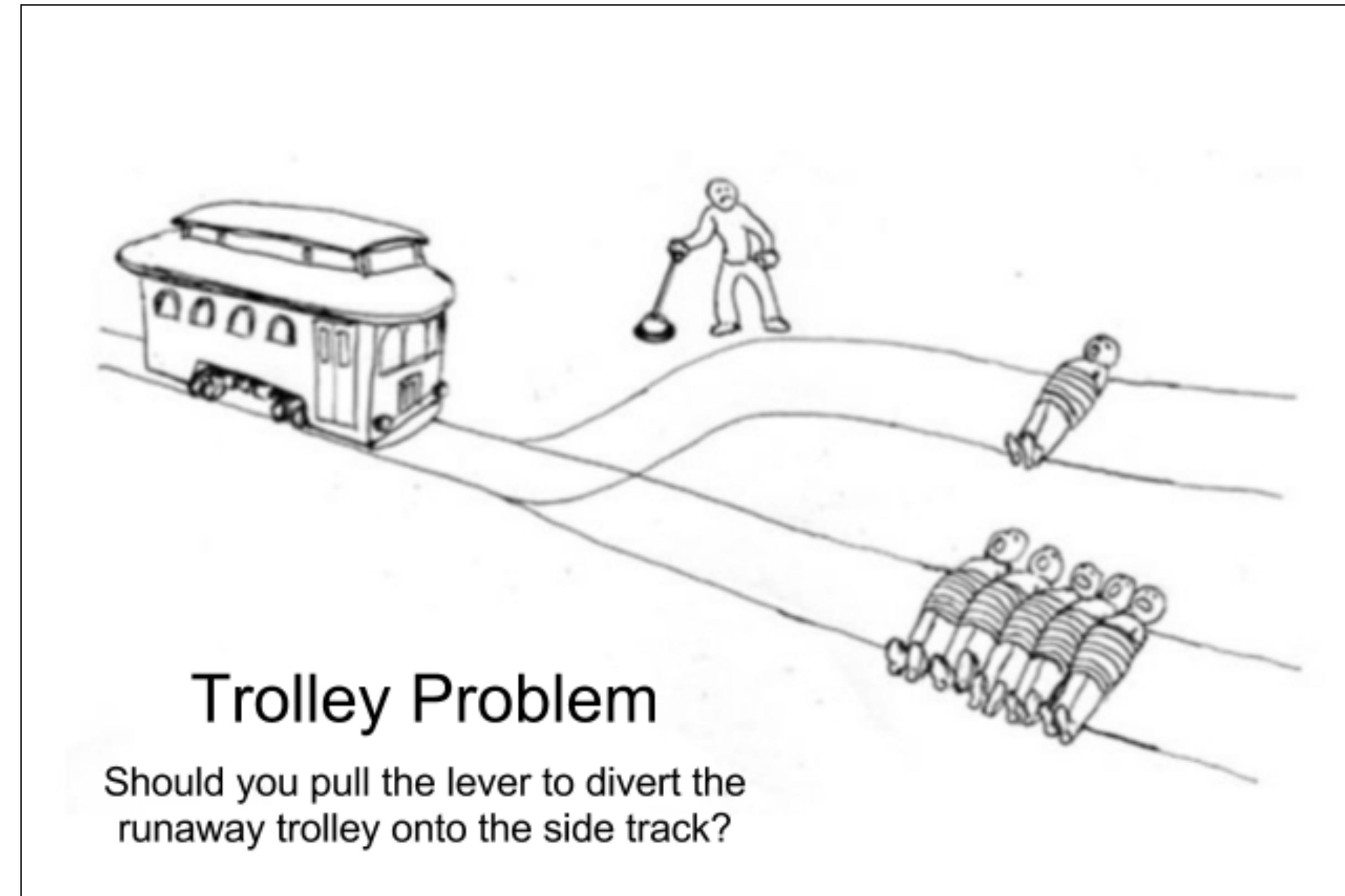
- Professional philosophers are just about evenly split between these

Normative ethics: deontology, consequentialism, or virtue ethics?

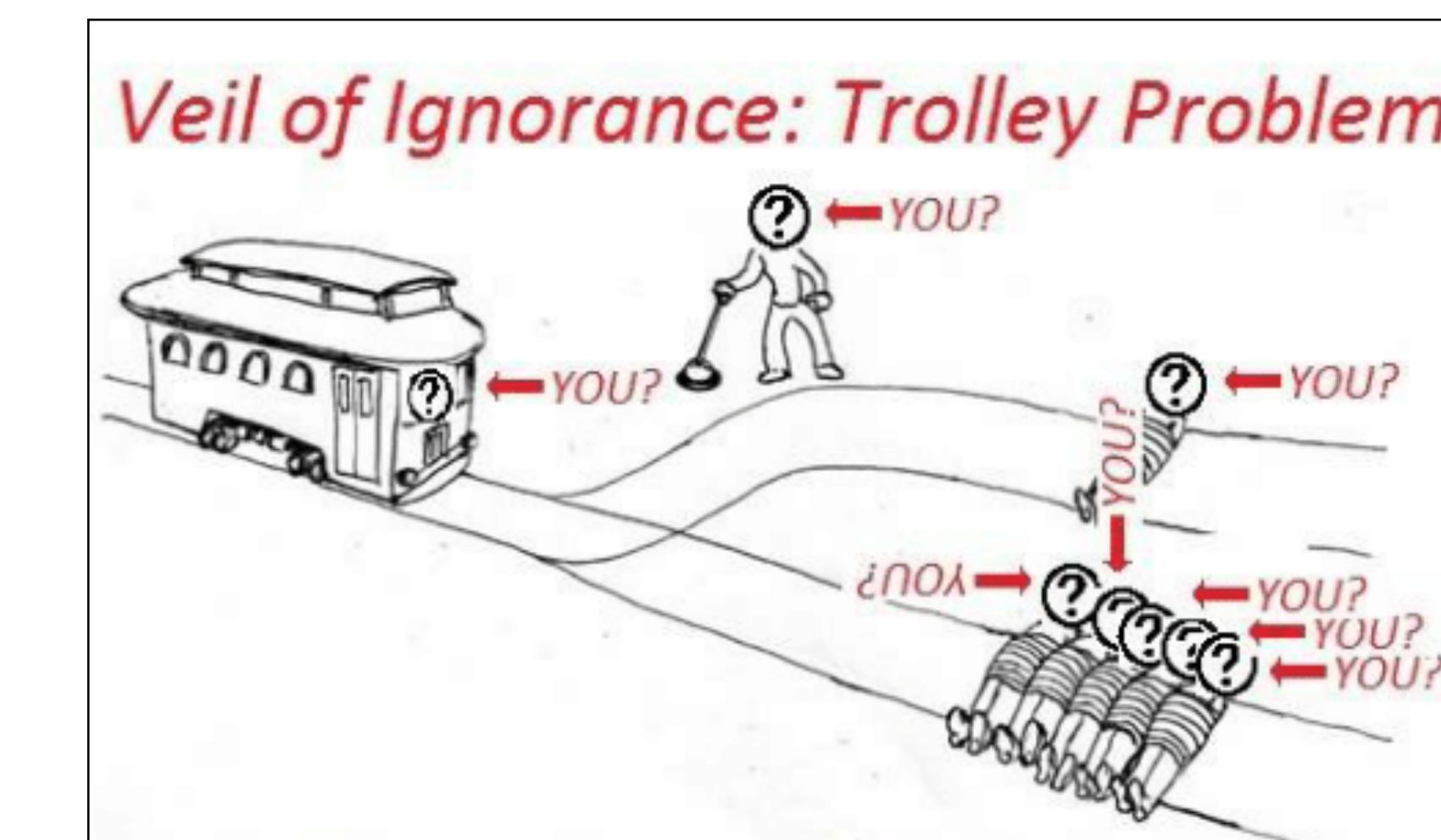
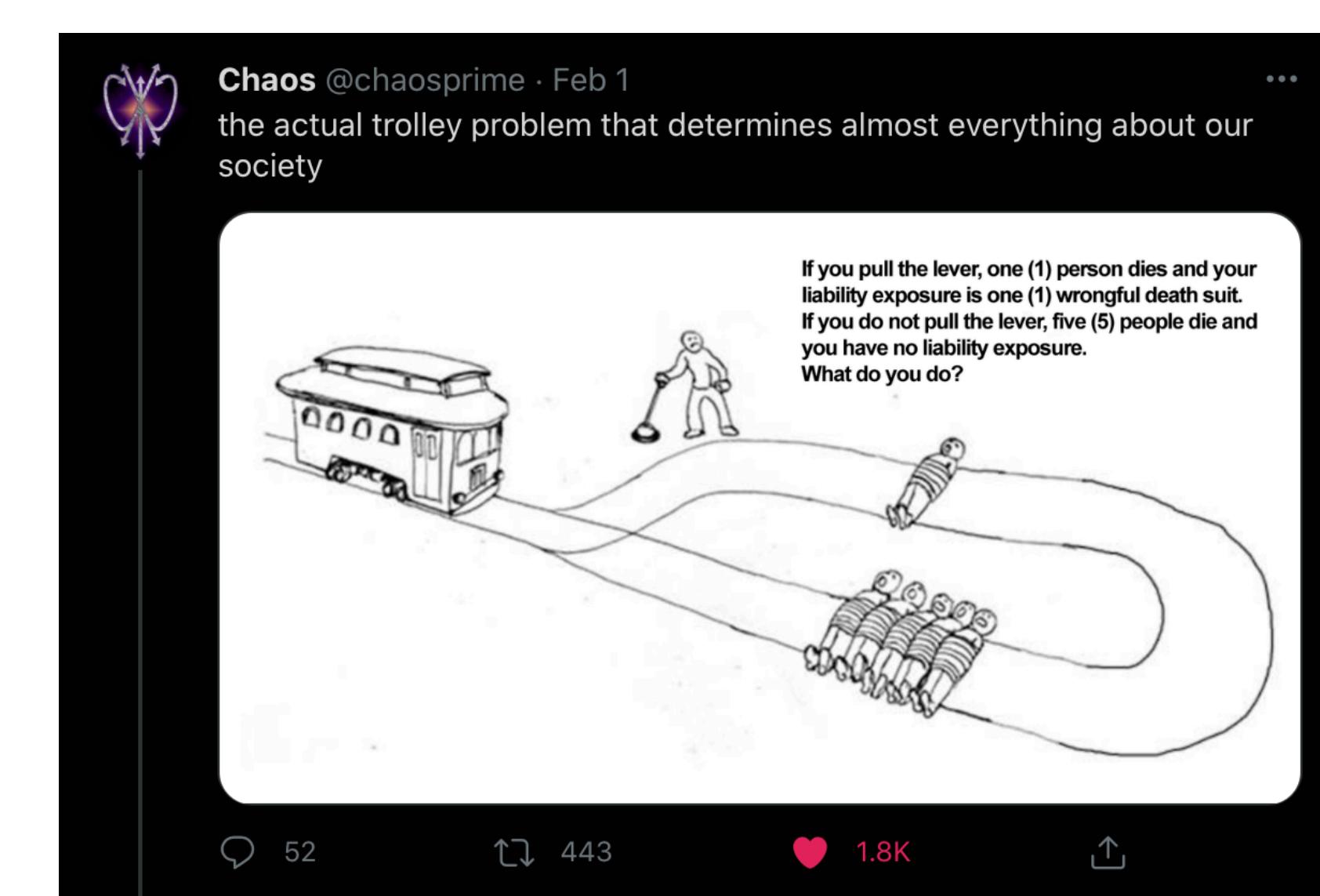
| | |
|---|-------------------|
| Other | 301 / 931 (32.3%) |
| Accept or lean toward: deontology | 241 / 931 (25.9%) |
| Accept or lean toward: consequentialism | 220 / 931 (23.6%) |
| Accept or lean toward: virtue ethics | 169 / 931 (18.2%) |

<https://philpapers.org/surveys/results.pl>

Intuition through "Trolley Problems"



Trolley Problems

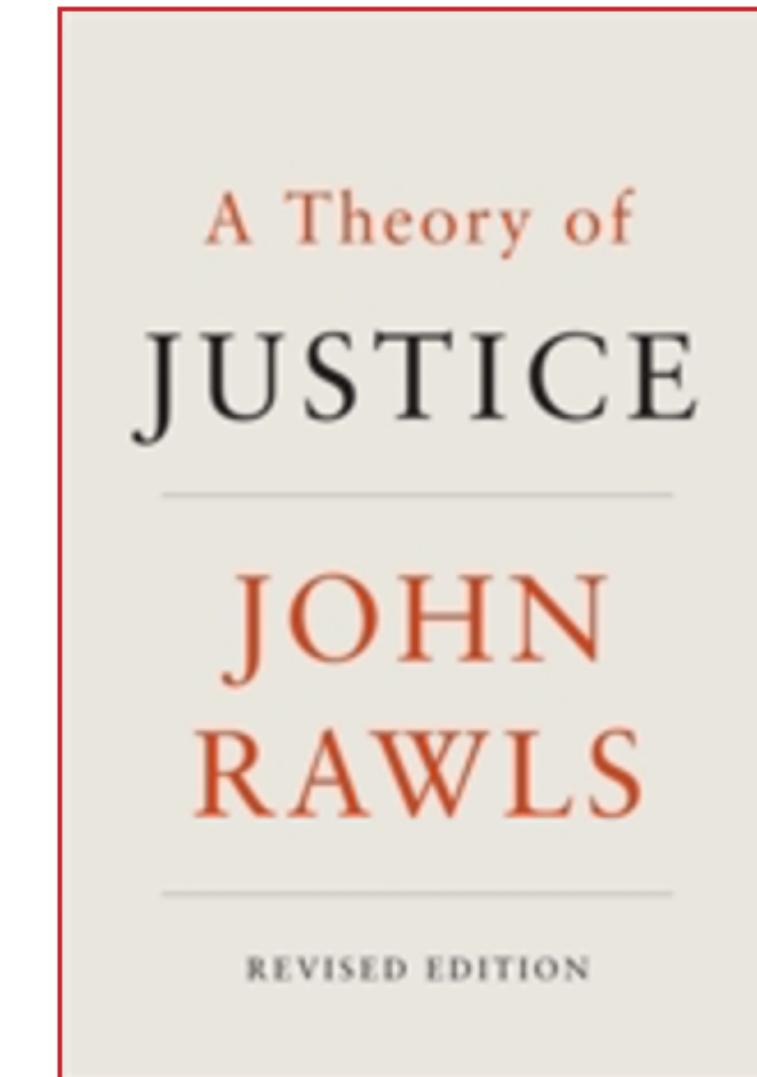


Ethical Theories

Rawls' theory of justice

(John Rawls, US, 1921 –)

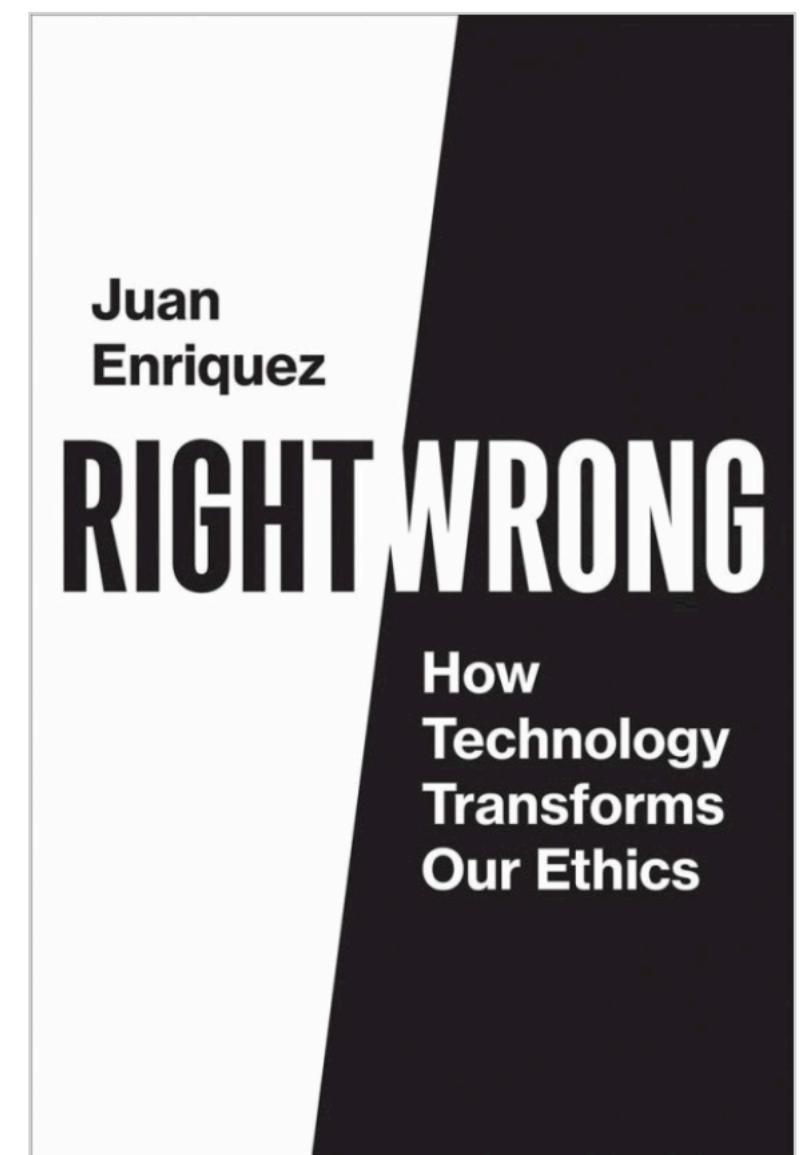
- Distributive justice
- Fair and equitable treatment of all based on a social contract
- Asks: what would you prescribe from under a 'veil of ignorance'?
- Any costs or consequences are secondary to this basic focus
- 'Do as you would be done by'



https://www.researchgate.net/figure/Summary-of-major-ethical-theories_tbl2_229658531

Ethics of Technology

- Ethics change with technological progress
- e.g. Industrial revolution
- e.g. Right to Internet access
- e.g. Birth control, surrogate pregnancy, embryo selection, artificial womb
- e.g. Lab-grown meat

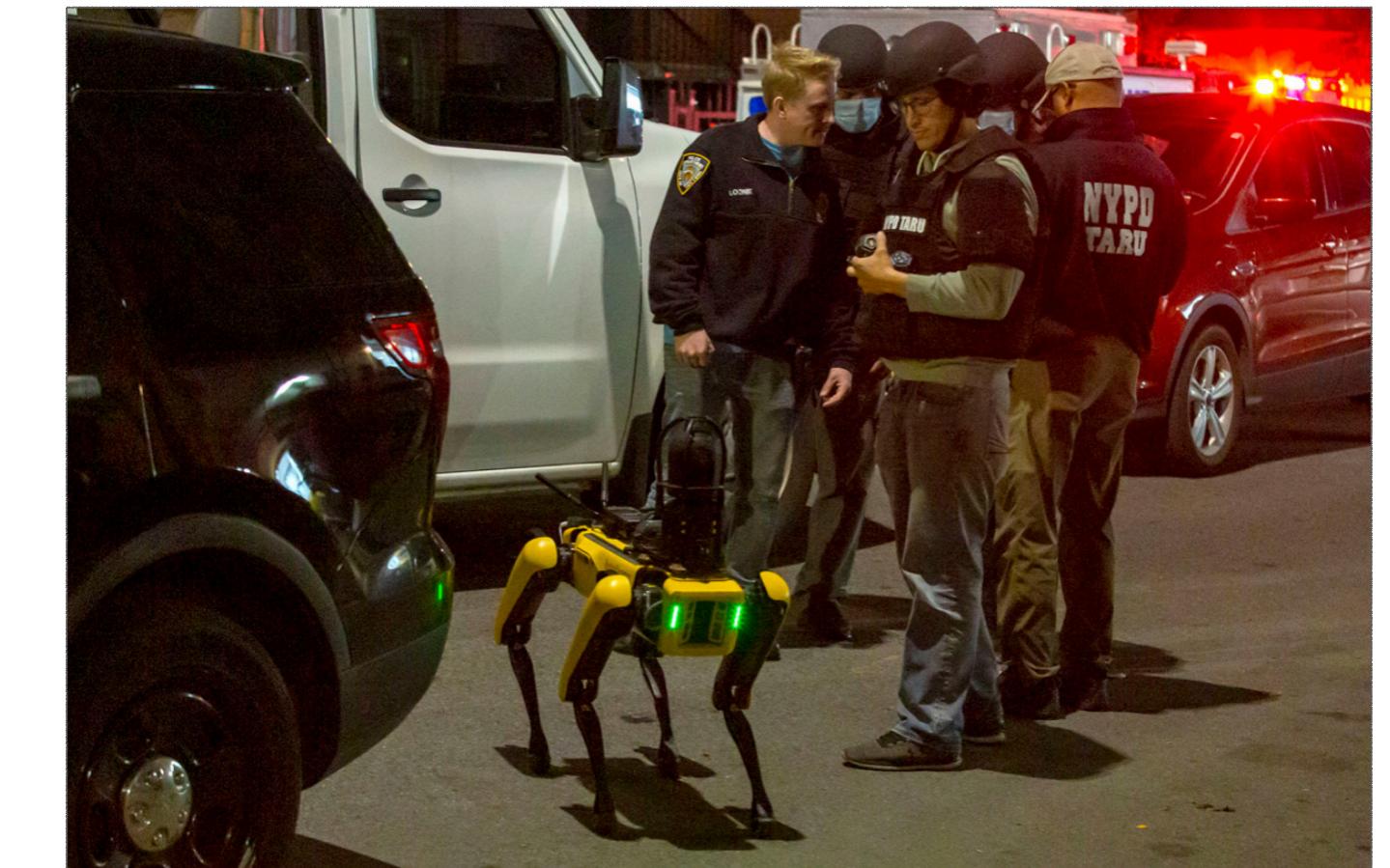
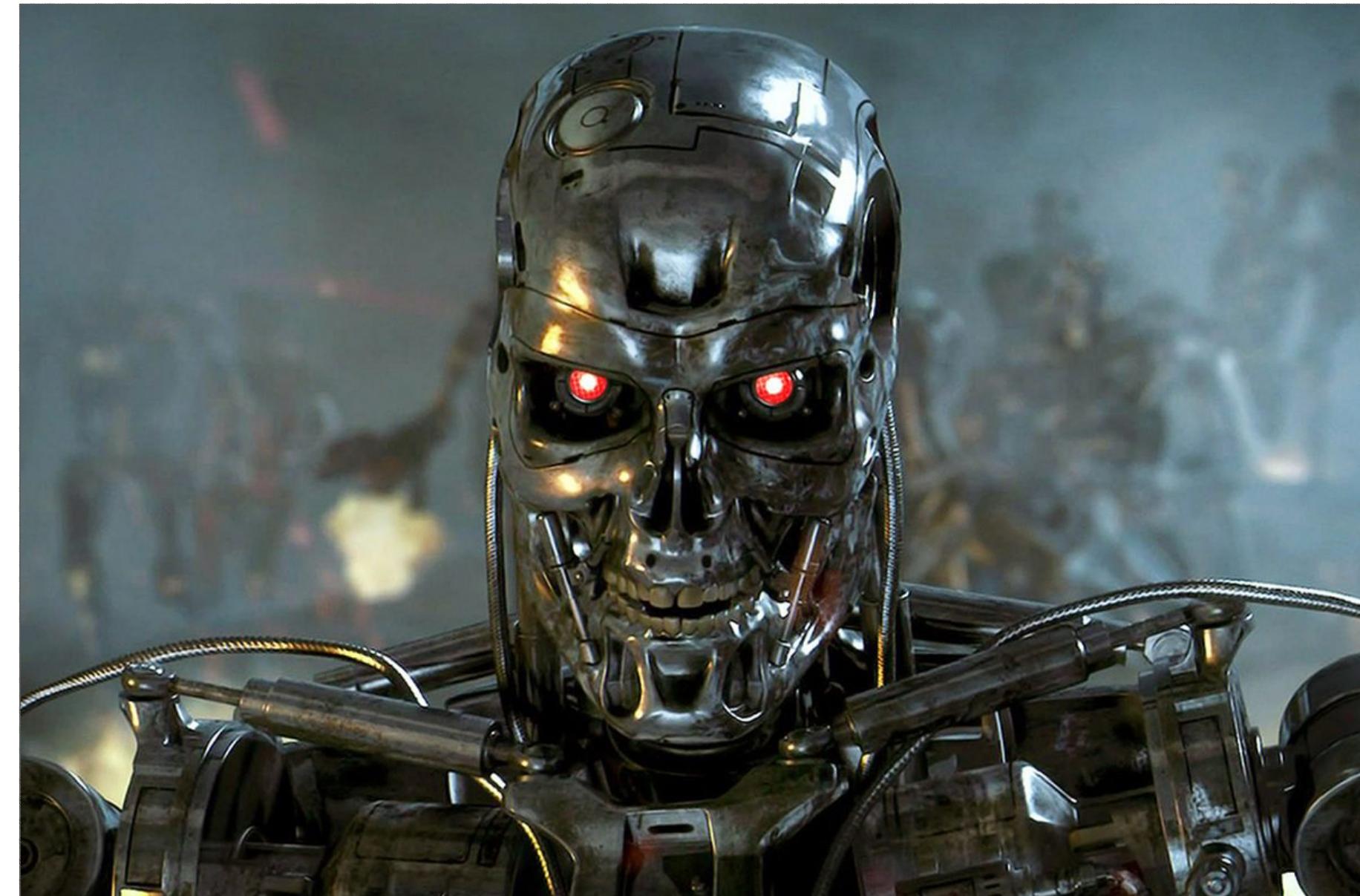


Questions?

Long-term problems

Autonomous Weapons

- Tempting to dismiss as far-fetched at this time
- But "the future is already here, just not evenly distributed"



Replacing human labor

RODNEY BROOKS *Robots, AI, and other stuff*

BLOG MIT ROBUST.AI



POST: MEGATREND: THE DEMOGRAPHIC INVERSION

APRIL 9, 2017 — ESSAYS

Megatrend: The Demographic Inversion
rodneybrooks.com/megatrend-the-demographic-inversion/

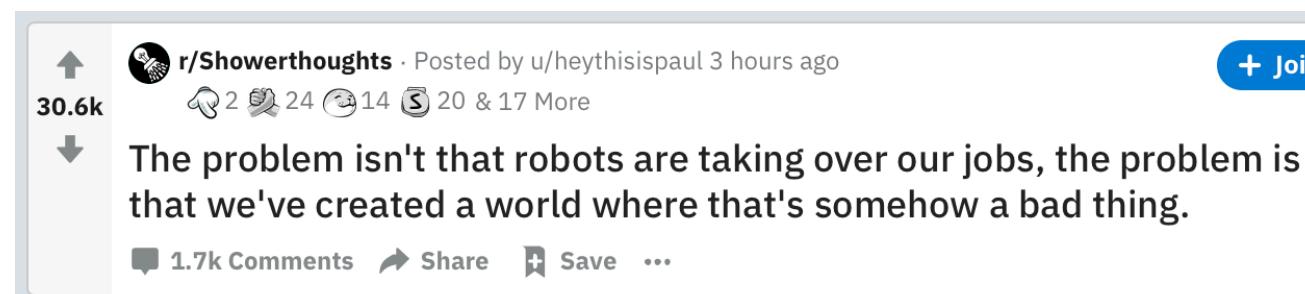
TIME

SUBSCRIBE

BUSINESS • COVID-19

Millions of Americans Have Lost Jobs in the Pandemic—And Robots and AI Are Replacing Them Faster Than Ever

- AI and robots replacing humans in existing jobs
 - Good and bad!
 - Interesting spin: AI controlling human labor

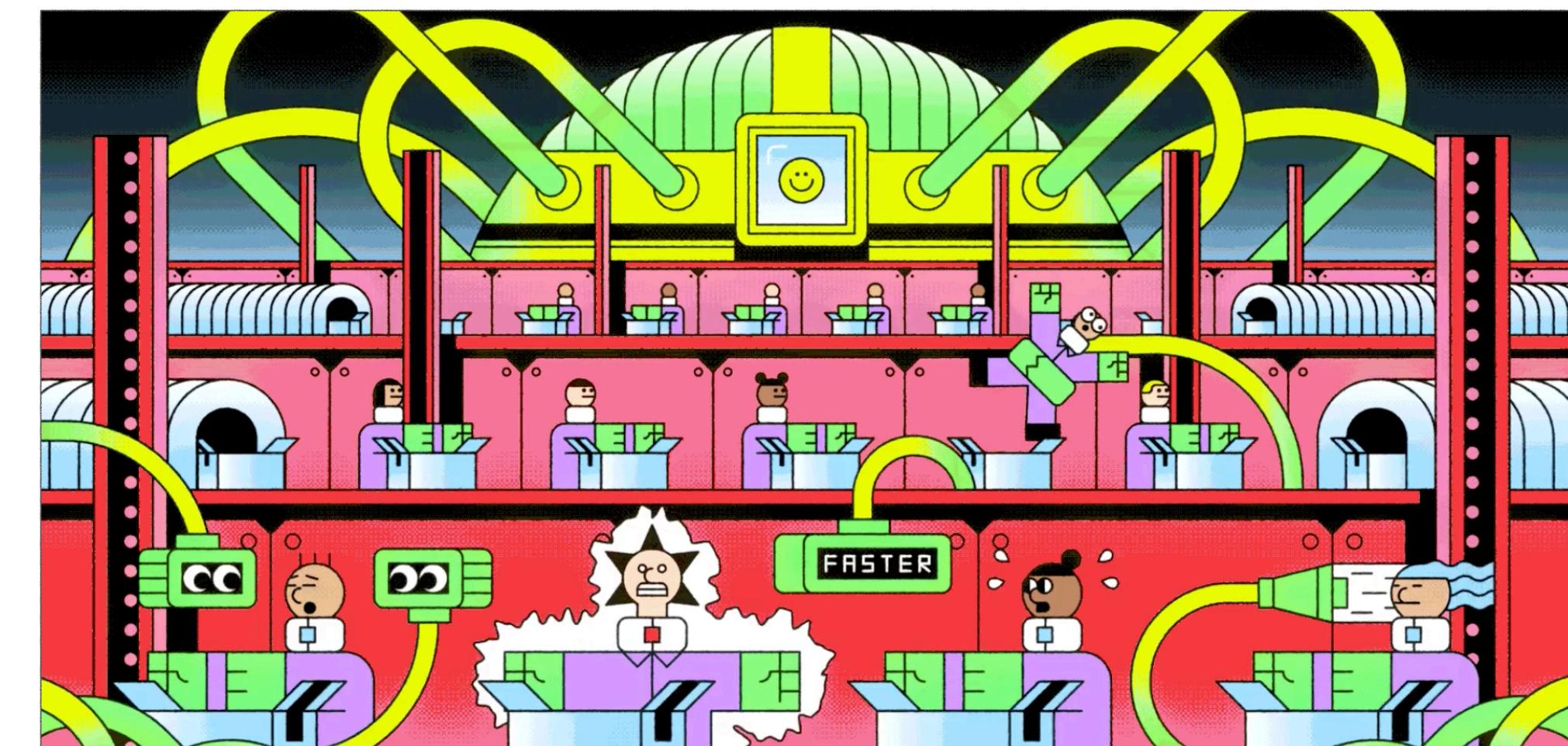


Recommended
short story →

Manna – Two Views of Humanity's Future – Chapter 1

by [Marshall Brain](#)

Depending on how you want to think about it, it was funny or inevitable or symbolic that the robotic takeover did not start at MIT, NASA, Microsoft or Ford. It started at a Burger-G restaurant in Cary, NC on May 17. It seemed like such a simple thing at the time, but May 17 marked a pivotal moment in human history.



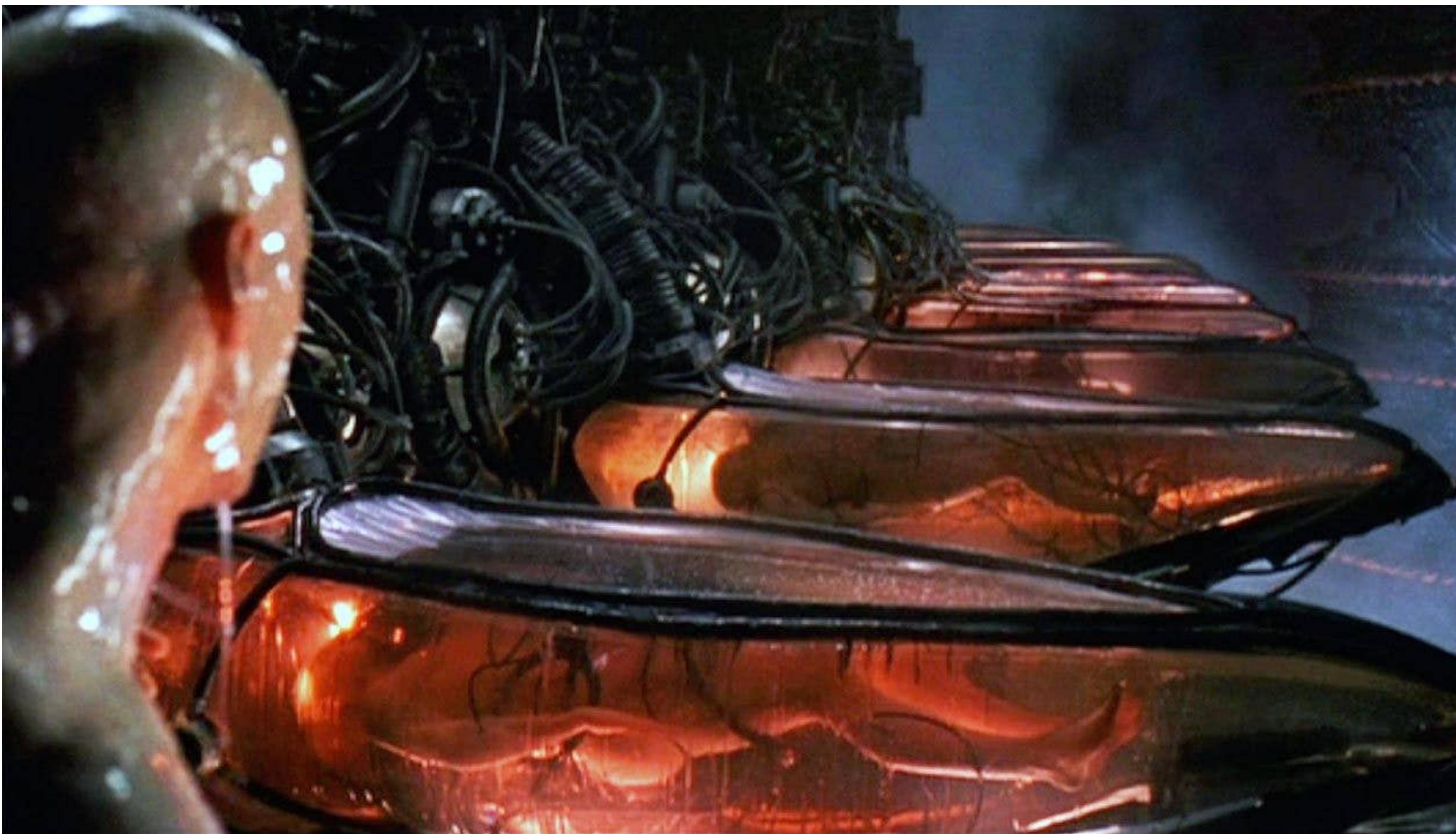
FEATURES

HOW HARD WILL THE ROBOTS MAKE US WORK?

In warehouses, call centers, and other sectors, intelligent machines are managing humans, and they're making work more stressful, grueling, and dangerous

By [Josh Dzieza](#) | [@joshdzieza](#) | Feb 27, 2020, 8:00am EST
Illustrations by [Joel Plosz](#)

Replacing humans entirely



**What's common in all these
long-term problems?**

Alignment

- "Paperclip maximizer": AGI given the goal of producing paperclips eventually turns every atom in space into paperclips
- Old lesson: establishing and communicating our goals and values is hard, and technology amplifies the difficulty



Guiding principle:
AI systems we build need to be
aligned with our goals and values

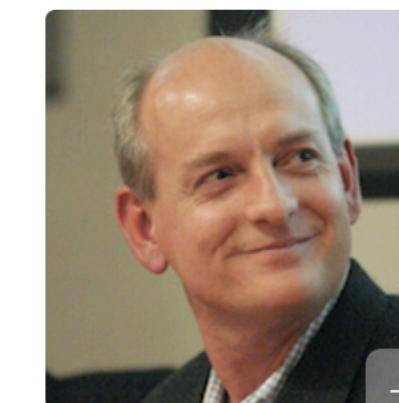
Questions?

Deep topic

- Active area of research (active at Berkeley)
- Useful lens for near-term problems, too



Faculty



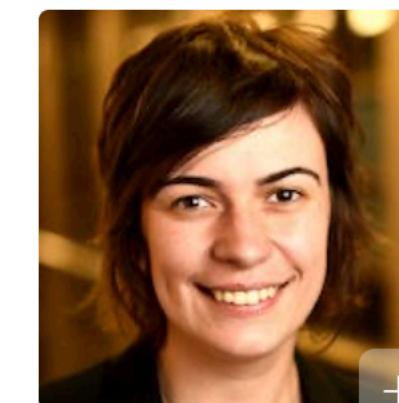
Stuart Russell

Professor of EECS at UC Berkeley
Holder of the Smith Zadeh Chair in Engineering
Author of *Artificial Intelligence: A Modern Approach*



Pieter Abbeel

Professor of EECS at UC Berkeley



Anca Dragan

Assistant Professor of EECS at UC Berkeley
Founder, InterACT Lab



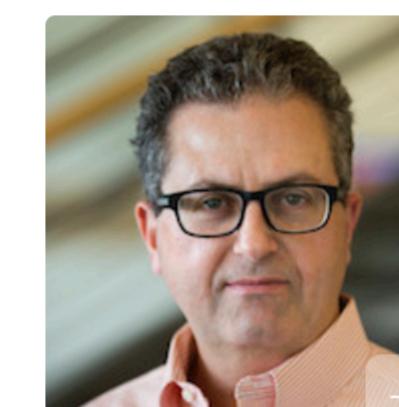
Bart Selman

Professor of Computer Science and Engineering at Cornell University



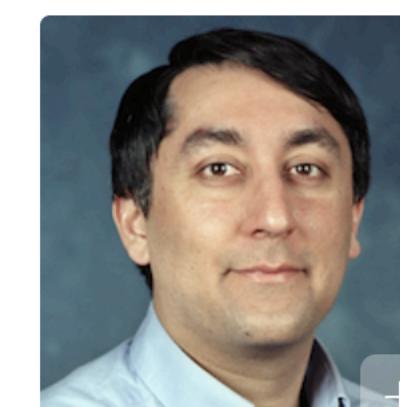
Joseph Halpern

Professor of Computer Science and Engineering at Cornell University



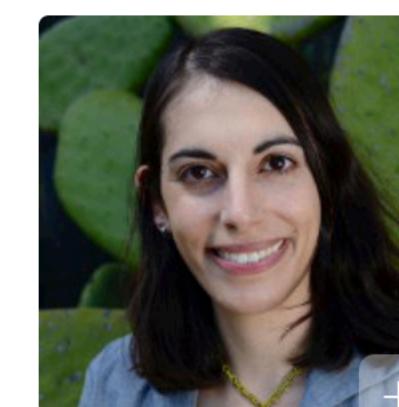
Michael Wellman

Professor of Computer Science and Engineering at the University of Michigan



Satinder Singh Baveja

Professor of Computer Science and Engineering at the University of Michigan
Director of the Artificial Intelligence Laboratory



Tania Lombrozo

Professor of Psychology at UC Berkeley



Tom Griffiths

Professor of Psychology and Cognitive Science at Princeton
Previously Director of the Computational Cognitive Science Lab
Previously Director of the

Near-term problems

Hiring

RETAIL OCTOBER 10, 2018 / 4:04 PM / UPDATED 2 YEARS AGO

Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

8 MIN READ

SAN FRANCISCO (Reuters) - Amazon.com Inc's [AMZN.O](#) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

Reuters

<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

Hiring

- ML model to predict hiring decision given a resume

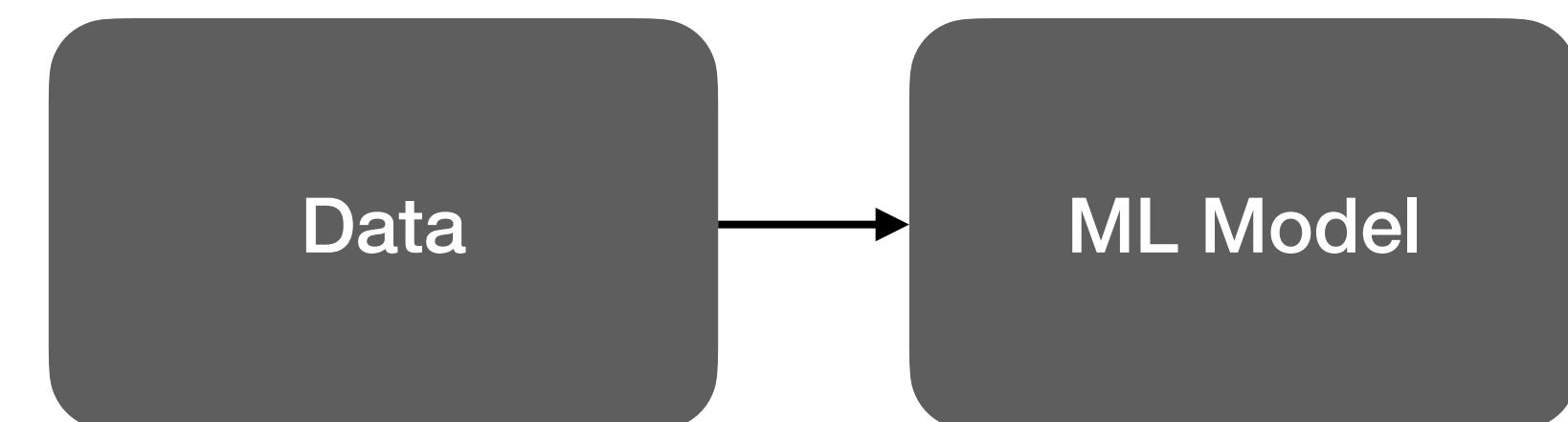


ML Model

Hiring

- ML model to predict hiring decision given a resume

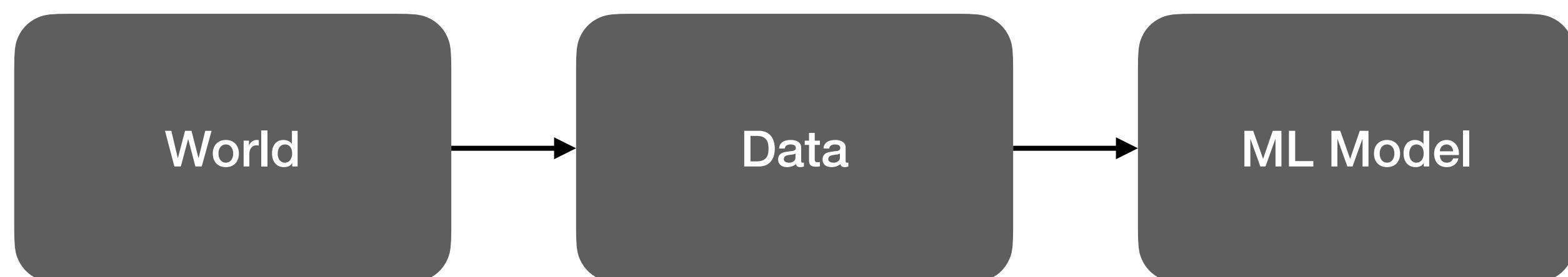
*Hiring decisions?
Or job performance?*



Hiring

- ML model to predict hiring decision given a resume

*So no matter what we pick,
our data is also biased*



Known to be biased in many ways:

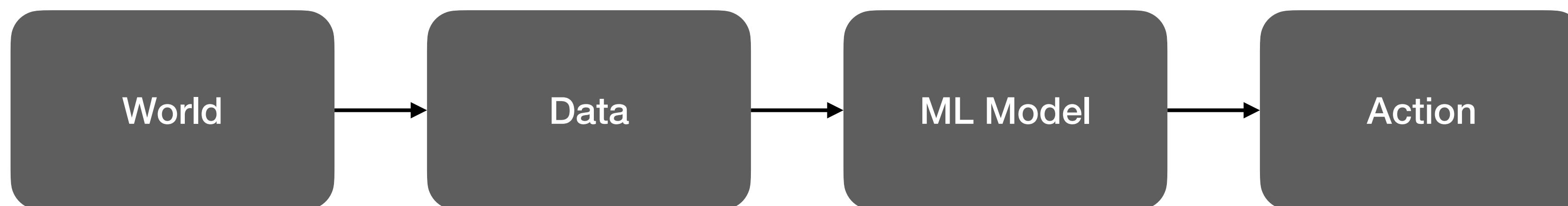
- hiring pipeline*
- hiring decisions*
- performance ratings*

*Therefore our
model is biased*

Hiring

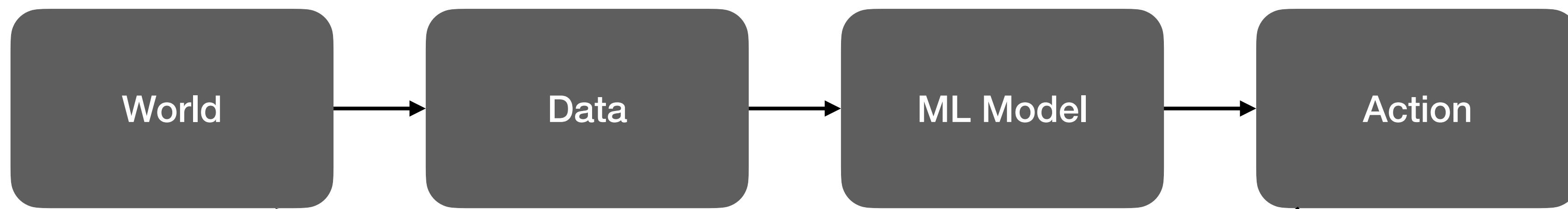
- ML model to predict hiring decision given a resume

*Sourcing? Double-checking human
decision? Hiring?*



Hiring

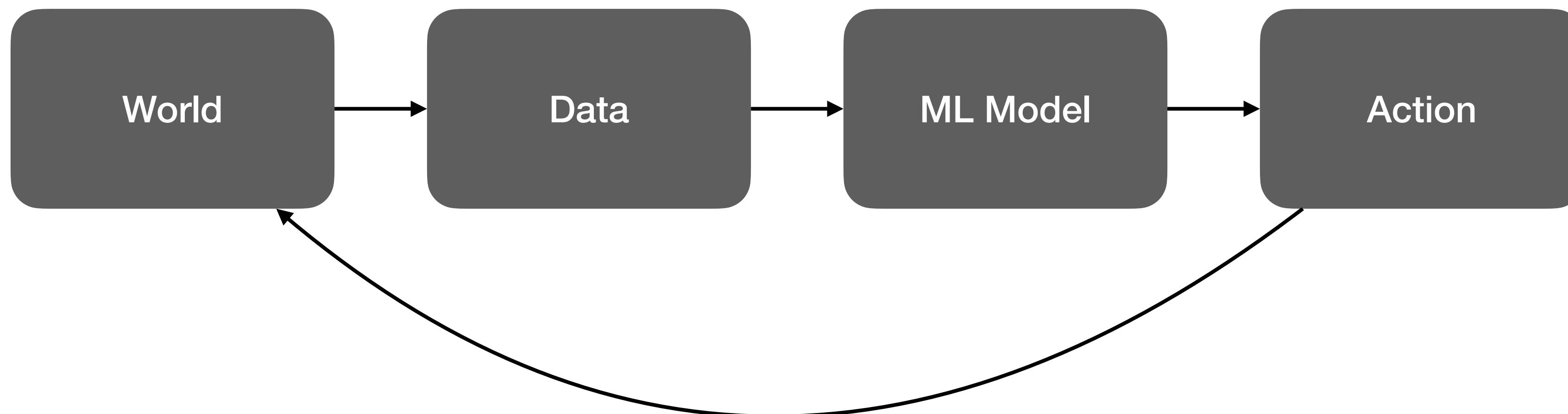
- ML model to predict hiring decision given a resume



*In many cases, amplifying
existing biases!*

Hiring

- ML model to predict hiring decision given a resume
- **Amplifying existing biases is not aligned with our goals and values!**



*In many cases, amplifying
existing biases!*

Questions?

Fairness

COMPAS

Correctional Offender Management Profiling for Alternative Sanctions

- Goal: predict recidivism, such that judges can consult 1-10 score in pre-trial sentencing decisions.
- Motivation: be less biased than humans
- Solution:
 - Gather data
 - Exclude protected class attributes (race, etc)
 - Ensure that our model's score corresponds to same probability of recidivism across all groups

And yet!

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

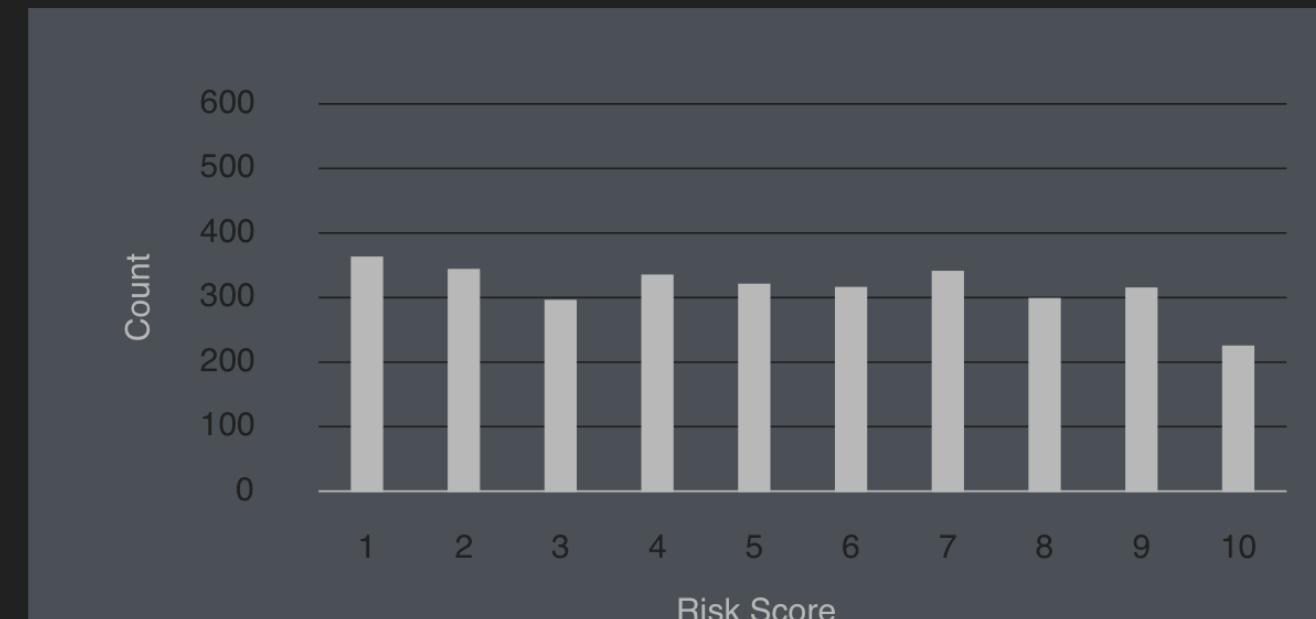
May 23, 2016

Prediction Fails Differently for Black Defendants

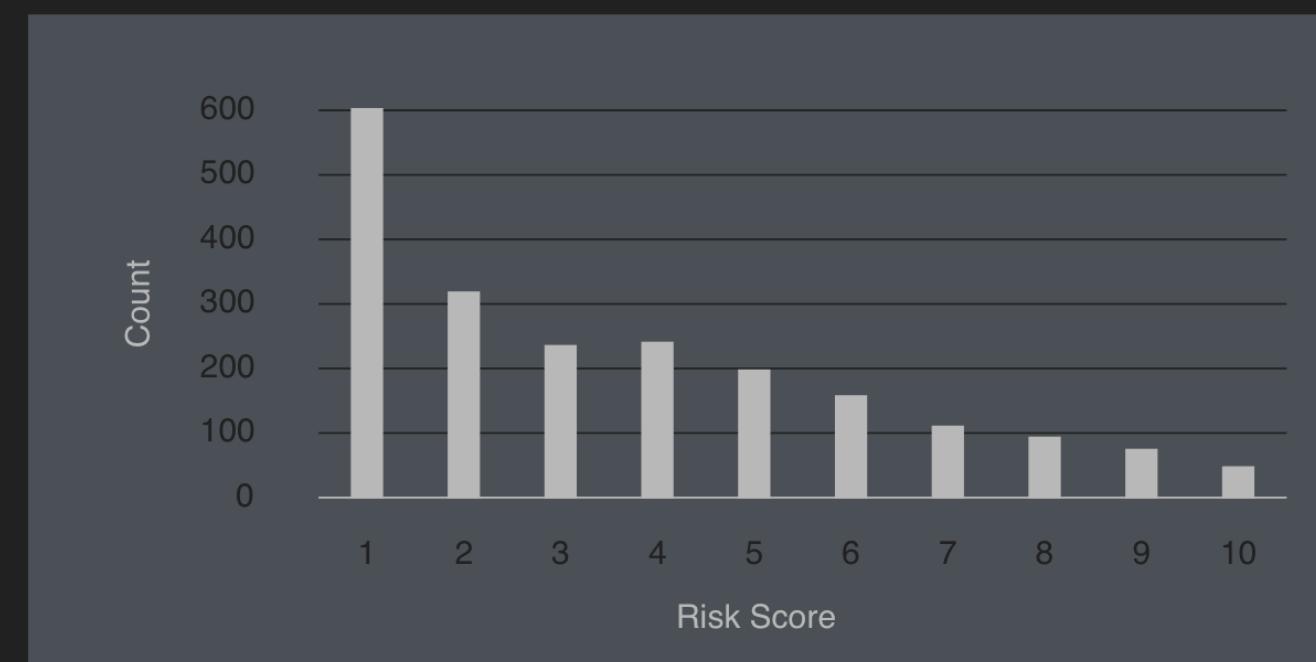
| | WHITE | AFRICAN AMERICAN |
|---|-------|------------------|
| Labeled Higher Risk, But Didn't Re-Offend | 23.5% | 44.9% |
| Labeled Lower Risk, Yet Did Re-Offend | 47.7% | 28.0% |

Overall, Northpointe's assessment tool correctly predicts recidivism 61 percent of the time. But blacks are almost twice as likely as whites to be labeled a higher risk but not actually re-offend. It makes the opposite mistake among whites: They are much more likely than blacks to be labeled lower risk but go on to commit other crimes. (Source: ProPublica analysis of data from Broward County, Fla.)

Black Defendants' Risk Scores

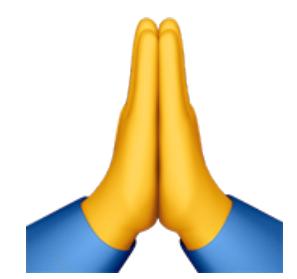


White Defendants' Risk Scores



These charts show that scores for white defendants were skewed toward lower-risk categories. Scores for black defendants were not. (Source: ProPublica analysis of data from Broward County, Fla.)

<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>



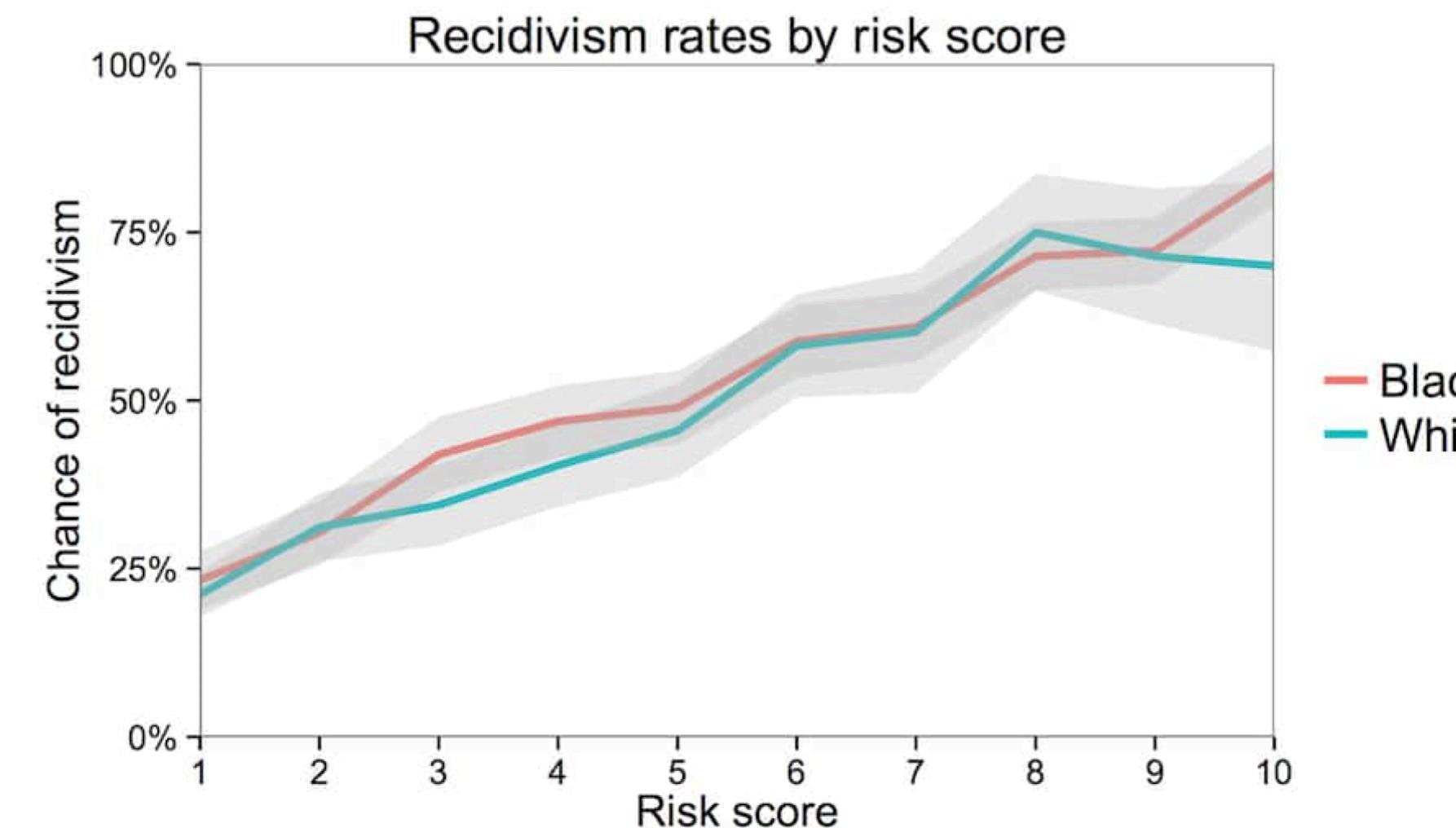
21 Fairness Definitions



Aravind Narayanan - 21 Fairness Definitions

Statistical Bias

- *Statistical bias:* Difference between estimator's expected value and the true value
- In this sense, COMPAS scores are not biased, w.r.t re-arrest ←
- But is it an adequate fairness criterion? ***Is it aligned with our values?***



Important caveat! We only have data for arrests, not crimes committed. There may well be bias in arrests.

<https://www.washingtonpost.com/news/monkey-cage/wp/2016/10/17/can-an-algorithm-be-racist-our-analysis-is-more-cautious-than-propublicas/>

Different perspectives on fairness

- What do different stakeholders want from the classifier?
- Decision-maker: of those I've labeled high risk, how many recidivated?
 - Predictive value: $TP / (TP + FP)$
- Defendant: what's the probability I'll be incorrectly classified as high risk?
 - False positive rate: $FP / (FP + TN)$
- Society: is the selected set demographically balanced?
 - Demographic parity

Problem setup: binary classification

| | | |
|---|----|----|
| | TN | FP |
| Did not recidivate | | |
| Recidivated | FN | TP |
| Labeled low-risk Labeled high-risk | | |

Looks oversimplified, but

- yields useful insights
- applicable to many contexts

Loans, hiring, insurance, etc.

- Loans: defaulted vs. didn't
- Hiring: succeeded at job vs. didn't

Group Fairness

- Do outcomes differ between groups (e.g. demographic), which we have no reason to believe are actually different?
- Motivation of the Pro-Publica article

Group fairness: impossibility theorem

if an instrument satisfies **predictive parity** ... but the prevalence differs between groups, the instrument cannot achieve **equal false positive [rates]** and **[equal] false negative rates** across those groups.

False Positive
False Negative

Prediction Fails Differently for Black Defendants

| | WHITE | AFRICAN AMERICAN |
|---|-------|------------------|
| Labeled Higher Risk, But Didn't Re-Offend | 23.5% | 44.9% |
| Labeled Lower Risk, Yet Did Re-Offend | 47.7% | 28.0% |

Overall, Northpointe's assessment tool correctly predicts recidivism 61 percent of the time. But blacks are almost twice as likely as whites to be labeled a higher risk but not actually re-offend. It makes the opposite mistake among whites: They are much more likely than blacks to be labeled lower risk but go on to commit other crimes. (Source: ProPublica analysis of data from Broward County, Fla.)

Chouldechova, *Fair Prediction with Disparate Impact: A study of bias in recidivism prediction instruments.*

Aravind Narayanan - 21 Fairness Definitions

No "correct" group fairness definition

Many group fairness metrics have natural motivations

Definitions

| Metric | Equalized under |
|-----------------------|---------------------|
| Selection probability | Demographic parity* |
| Pos. predictive value | Predictive parity |
| Neg. predictive value | |
| False positive rate | Error rate balance |
| False negative rate | Error rate balance |
| Accuracy | Accuracy equity |

Chouldechova paper

*aka disparate impact and many variants

Different metrics matter to different stakeholders.
There is no "right" definition.

If these metrics are equal for 2 groups, some trivial algebra shows that the prevalence is also the same between the groups

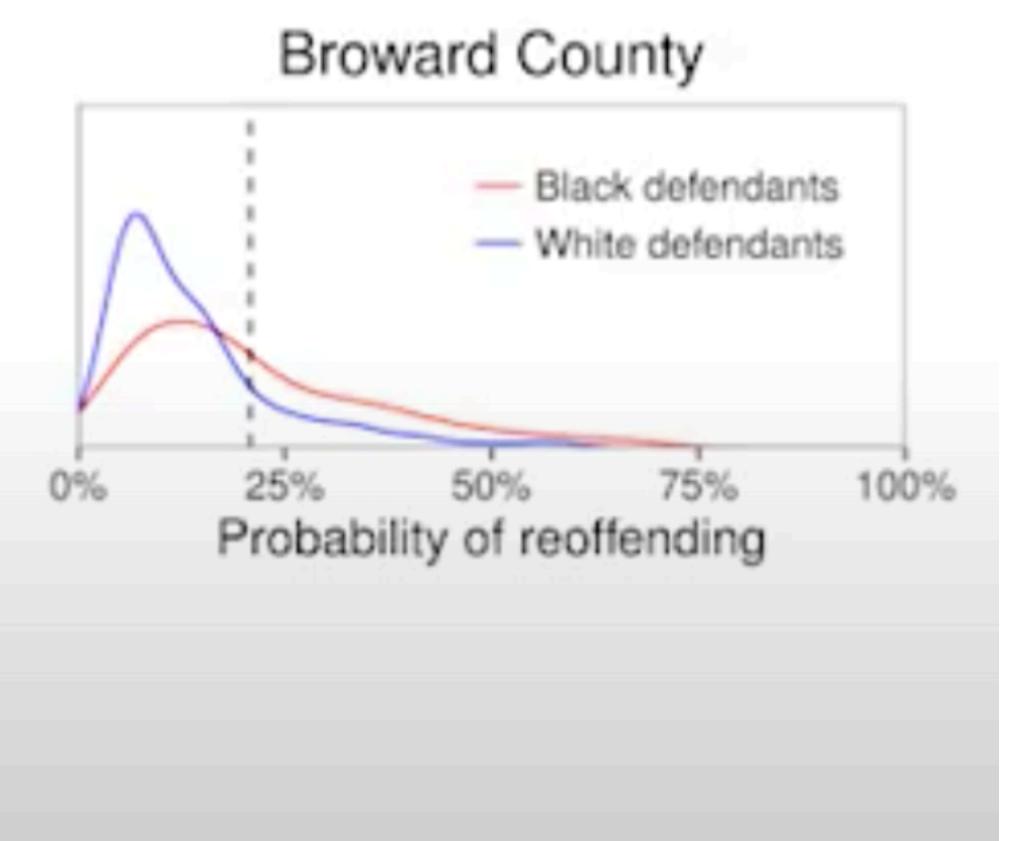
But there's nothing special about these 3! We can pick any 3

Aravind Narayanan - 21 Fairness Definitions

It gets worse...

- Okay, then let's just pick most important **two** metrics (FPR and FNR), and allow the model to use protected class attributes.
 - Now we fail individual fairness ->
- Fine, let's just pick **one!** (e.g. equal FPR)
 - We still sacrifice some *utility* (e.g. public safety, or number of defendants released)

Generally impossible to pick a single threshold for all groups that equalizes both FPR & FNR



*Assuming the scores are calibrated

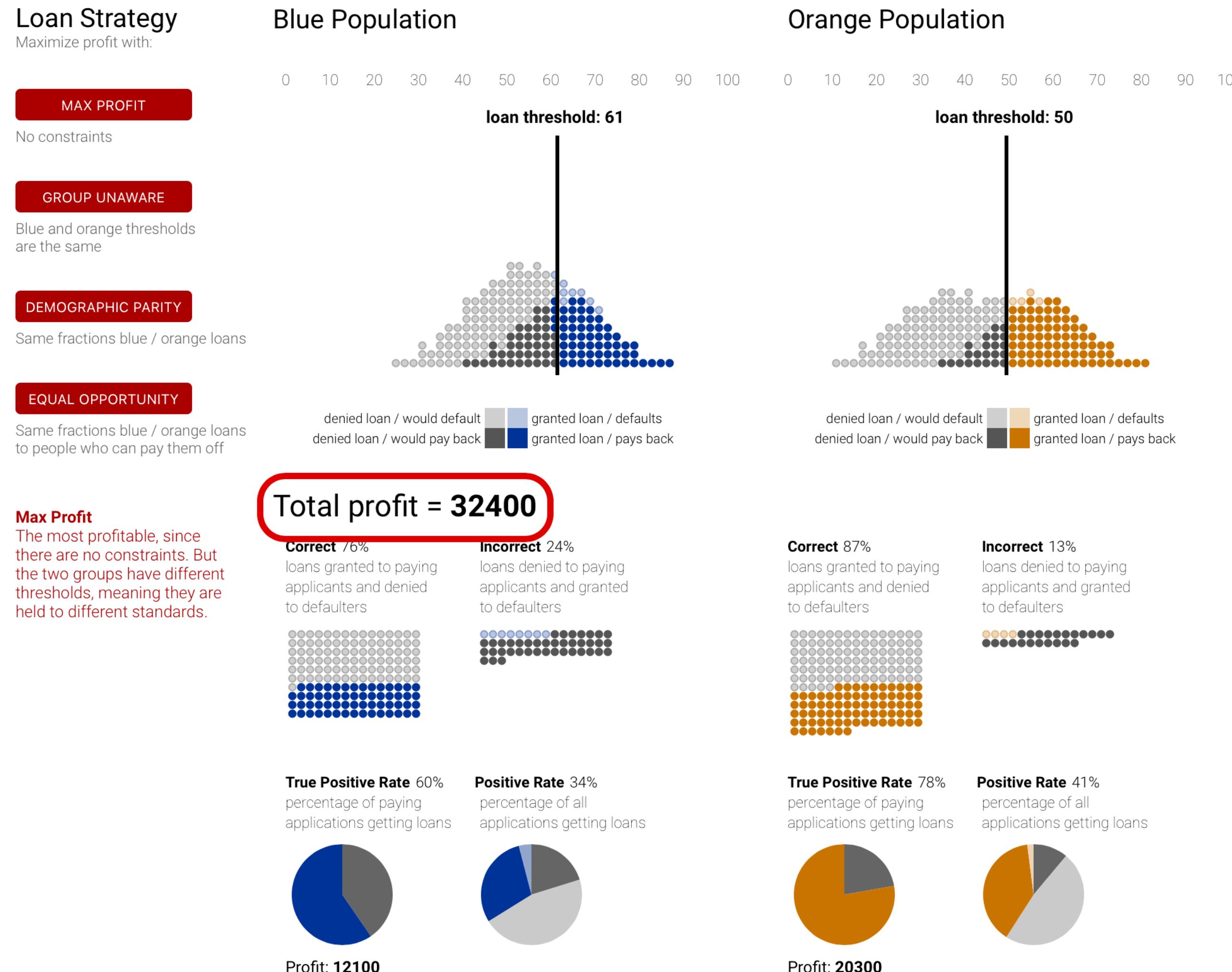
Corbett-Davies et al., *Algorithmic Decision Making and the Cost of Fairness*

Aravind Narayanan - 21 Fairness Definitions

Interactive Example

Simulating loan decisions for different groups

Drag the black threshold bars left or right to change the cut-offs for loans.
Click on different preset loan strategies.

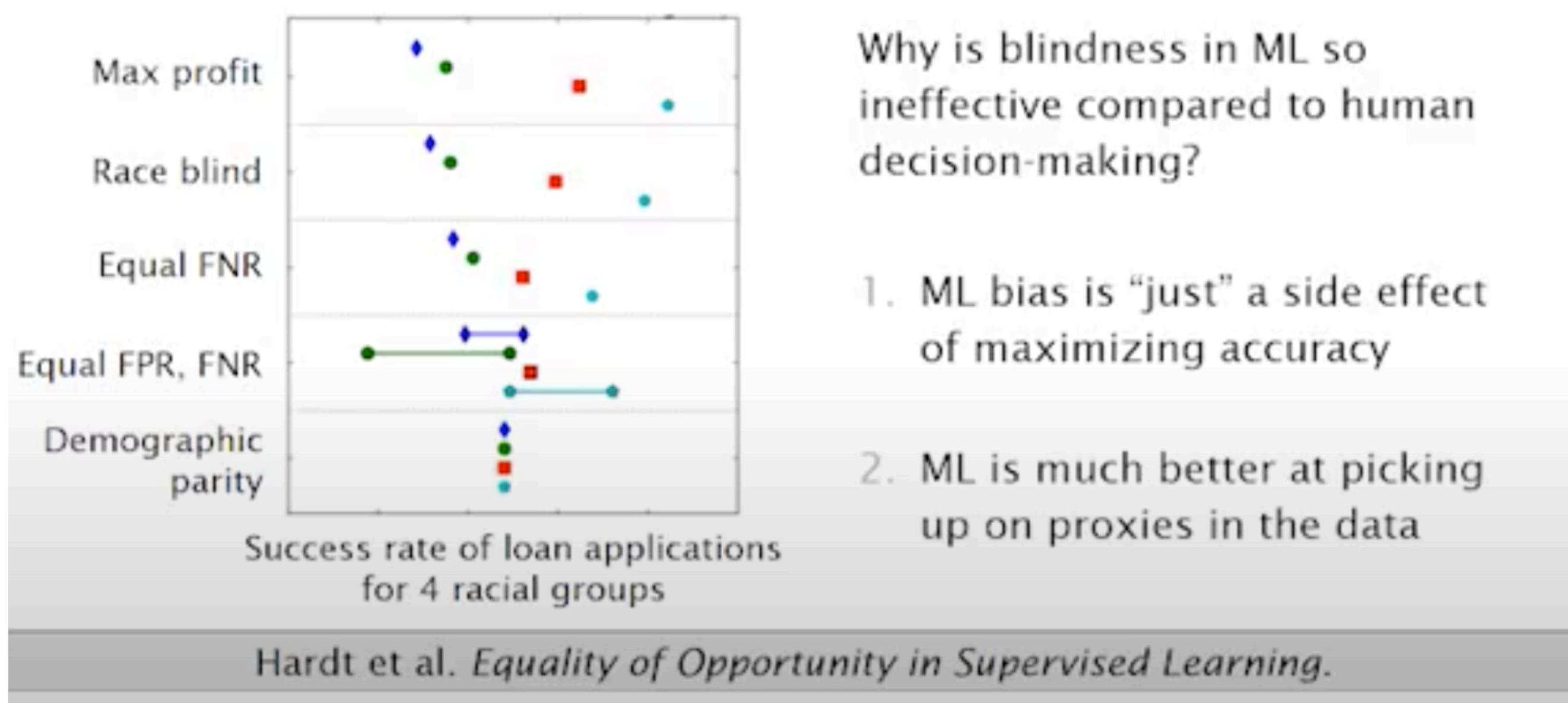


<https://research.google.com/bigpicture/attacking-discrimination-in-ml/>

By the way...

Ineffectiveness of blindness,
quantified

Provocation



Aravind Narayanan - 21 Fairness Definitions

Trade-offs

Provocation

How to navigate these tradeoffs?

- Between different measures of group fairness
- Between group fairness and individual fairness
- Between fairness and utility

Tension between disparate treatment
and disparate impact

Well known; see, e.g., Ricci vs. DeStefano

Finding creative case-by-case workarounds doesn't
"scale" for algorithmic decision making

New complexities: what about a system that uses the
protected attribute during training but not during
prediction?

Lipton et al. *Does mitigating ML's disparate impact require disparate treatment?*

Aravind Narayanan - 21 Fairness Definitions

Provocation

Much of this applies to human
decision making

Tensions between prediction desiderata
don't depend on who is doing the predicting

Sometimes in non-obvious ways; for example:
police search of a vehicle can be seen as a "prediction"
of the presence of contraband

Seeing the water

 **Moritz Hardt**
@mrtz ...

Apropos of current discussions, here's an example I found helpful in understanding why opting for prediction as a solution concept on its own (regardless of data and modeling choices) is already a consequential political act that deprioritizes alternatives.

Failure to appear in court

One approach: Predict failure to appear, jail if risk is high.

Alternative: Recognize that people fail to appear in court due to lack of child care and transportation, work schedules, or too many court appointments. Implement steps to mitigate these issues.

Alternative is part of the Harris County Lawsuit settlement: "*require Harris County to provide free child care at courthouses, develop a two-way communication system between courts and defendants, give cell phones to poor defendants and pay for public transit or ride share services for defendants without access to transportation to court.*" (Source: [Houston Chronicle, April 2019](#))



3:31 PM · Jun 23, 2020 · Twitter Web App

Equality



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

<https://www.mobilizegreen.org/blog/2018/9/30/environmental-equity-vs-environmental-justice-whats-the-difference>

Equality



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

Equity



Everyone gets the supports they need (this is the concept of "affirmative action"), thus producing equity.

<https://www.mobilizegreen.org/blog/2018/9/30/environmental-equity-vs-environmental-justice-whats-the-difference>

Equality



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

Equity



Everyone gets the supports they need (this is the concept of "affirmative action"), thus producing equity.

Justice



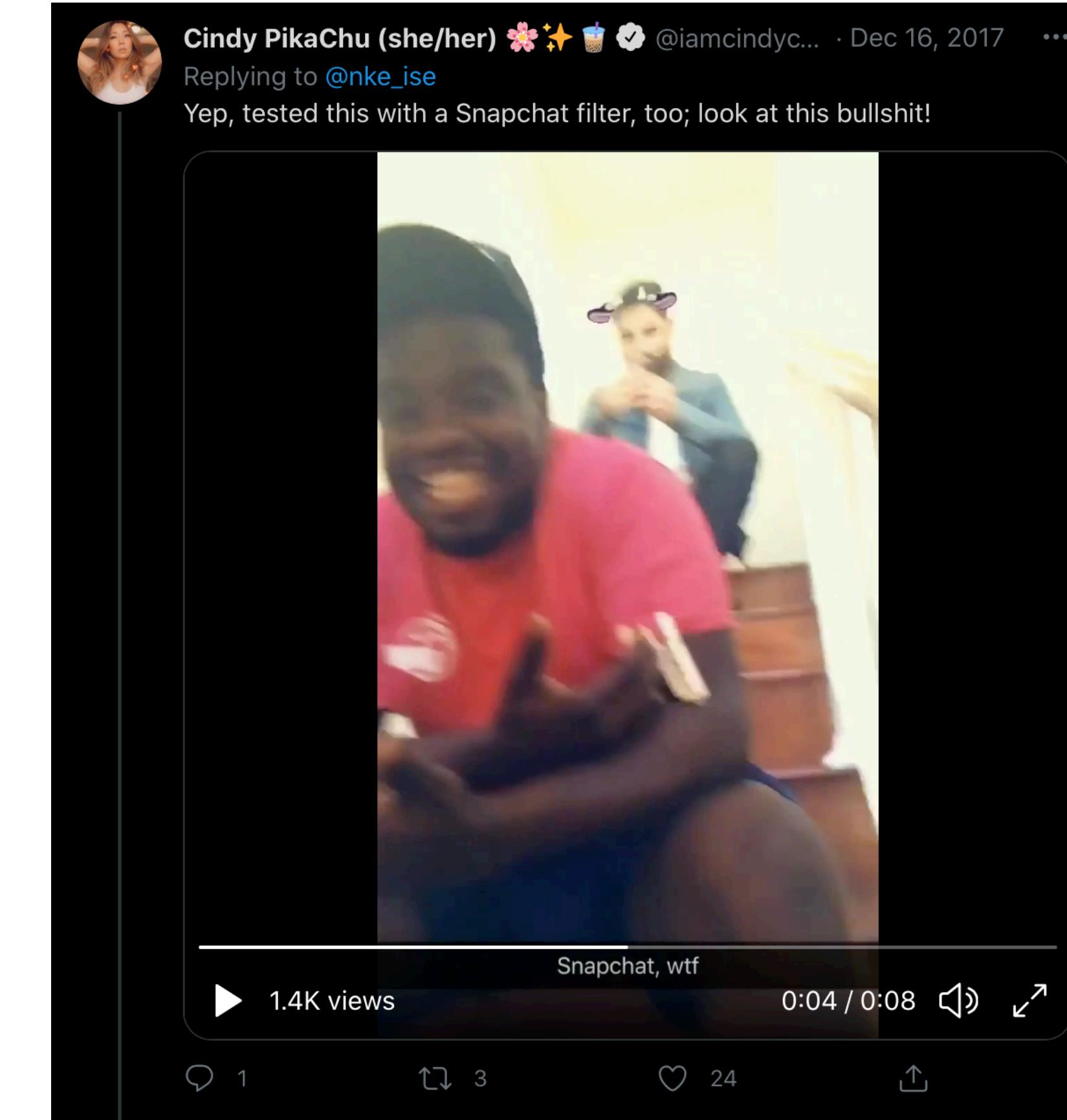
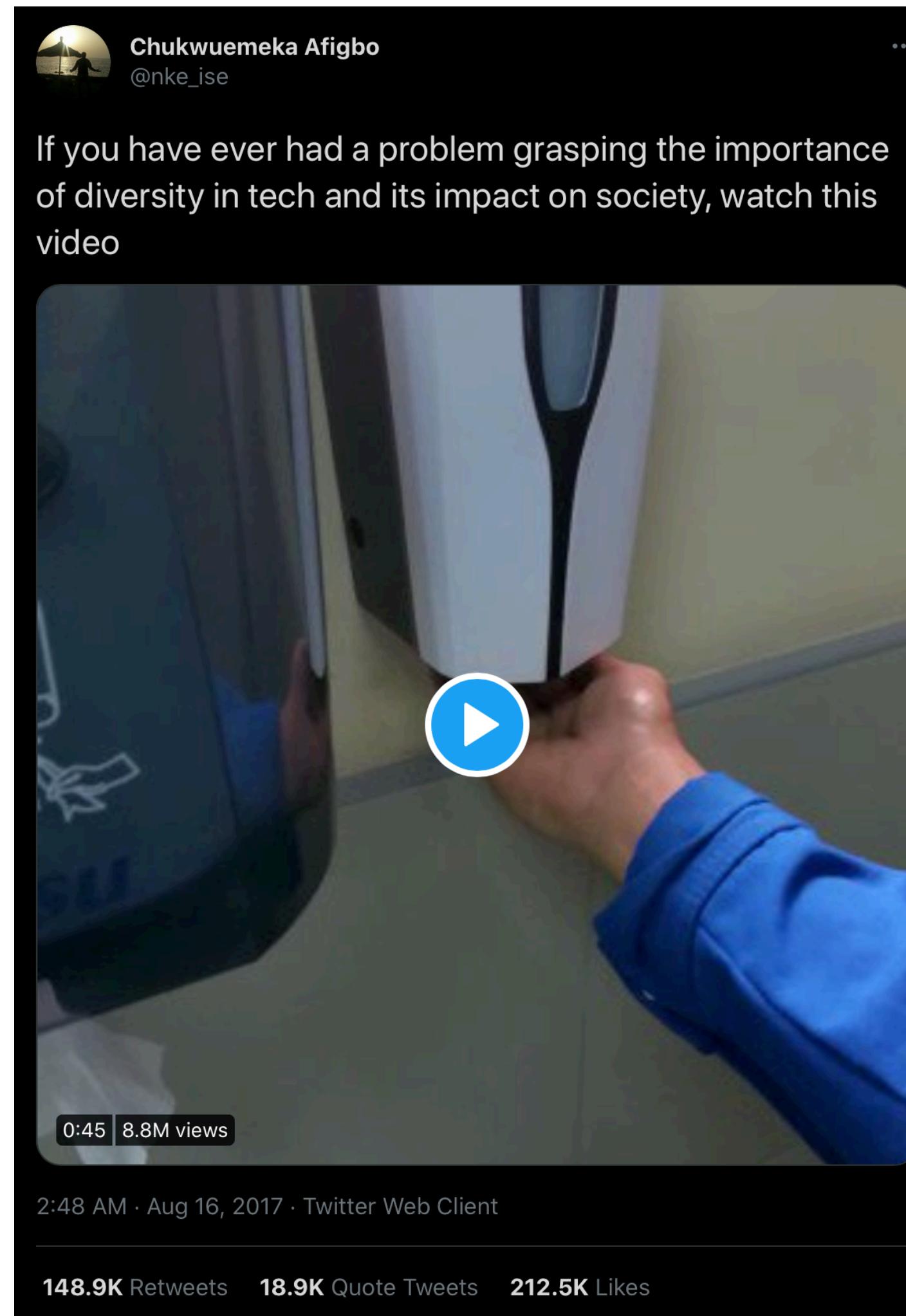
All 3 can see the game without supports or accommodations because **the cause(s) of the inequity was addressed**. The systemic barrier has been removed.

<https://www.mobilizegreen.org/blog/2018/9/30/environmental-equity-vs-environmental-justice-whats-the-difference>

Questions?

Representation

Representation



https://twitter.com/nke_isé/status/897756900753891328

Not a new problem, sadly

How Kodak's Shirley Cards Set Photography's Skin-Tone Standard

November 13, 2014 · 3:45 AM ET

Heard on [Morning Edition](#)



MANDALIT DEL BARCO



6-Minute Listen

+ PLAYLIST



Jersson Garcia works at Richard Photo Lab in Hollywood. He's 31 years old, and he's got a total crush on Shirley.

"Beautiful skin tones, beautiful eyes, great hair," he sighs.
"She's gorgeous."

Garcia is holding a 4-by-6-inch photo of an ivory-faced brunette wearing a lacy, white, off-the-shoulders top. She has red lipstick and silver earrings, and the photo appears to have been taken sometime in the 1970s or '80s.

For many years, this "Shirley" card — named for the original model, who was an employee of Kodak — was used by photo labs to calibrate skin tones, shadows and light during the printing process.



For decades, Kodak's Shirley cards, like this one, featured only white models.

Kodak

...or a problem in just our field

Lack of females in drug dose trials leads to overmedicated women

Gender gap leaves women experiencing adverse drug reactions nearly twice as often as men, study shows

Date: August 12, 2020

Source: University of California - Berkeley

Summary: Women are more likely than men to suffer adverse side effects of medications because drug dosages have historically been based on clinical trials conducted on men, suggests new research.

Recent improvement!

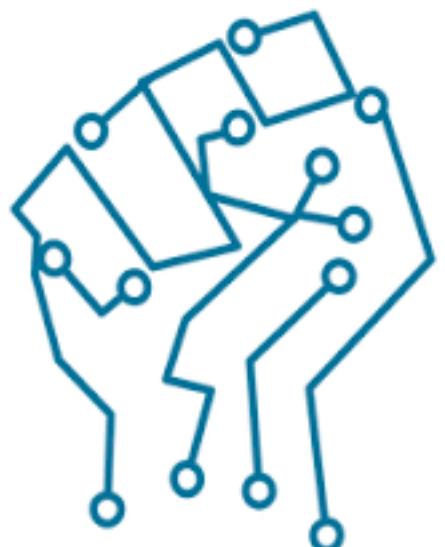
The screenshot shows the NPR Shots website. At the top, there's a dark header with the npr logo, a red 'DONATE' button, a 'Play Live Radio' button with a play icon, and links for 'HOURLY NEWS', 'LISTEN LIVE', and 'PLAYLIST'. Below the header, the 'Shots' logo is displayed next to a stylized syringe icon, with the text 'HEALTH NEWS FROM NPR'. A 'TREATMENTS' section features a headline: 'As COVID-19 Vaccine Trials Move At Warp Speed, Recruiting Black Volunteers Takes Time'. The date 'September 11, 2020 · 3:03 PM ET' and author 'BLAKE FARMER' are listed below the headline. A large blue button at the bottom left says '4-Minute Listen' with a play icon. To the right of the button are icons for '+ PLAYLIST', download ('download'), and share ('share').

Large part of the solution

The New York Times

Who Is Making Sure the A.I. Machines Aren't Racist?

When Google forced out two well-known artificial intelligence experts, a long-simmering research controversy burst into the open.



Hundreds of people gathered for the first lecture at what had become the world's most important conference on artificial intelligence — row after row of faces. Some were East Asian, a few were Indian, and a few were women. But the vast majority were white men. More than 5,500 people attended the meeting, five years ago in Barcelona, Spain.

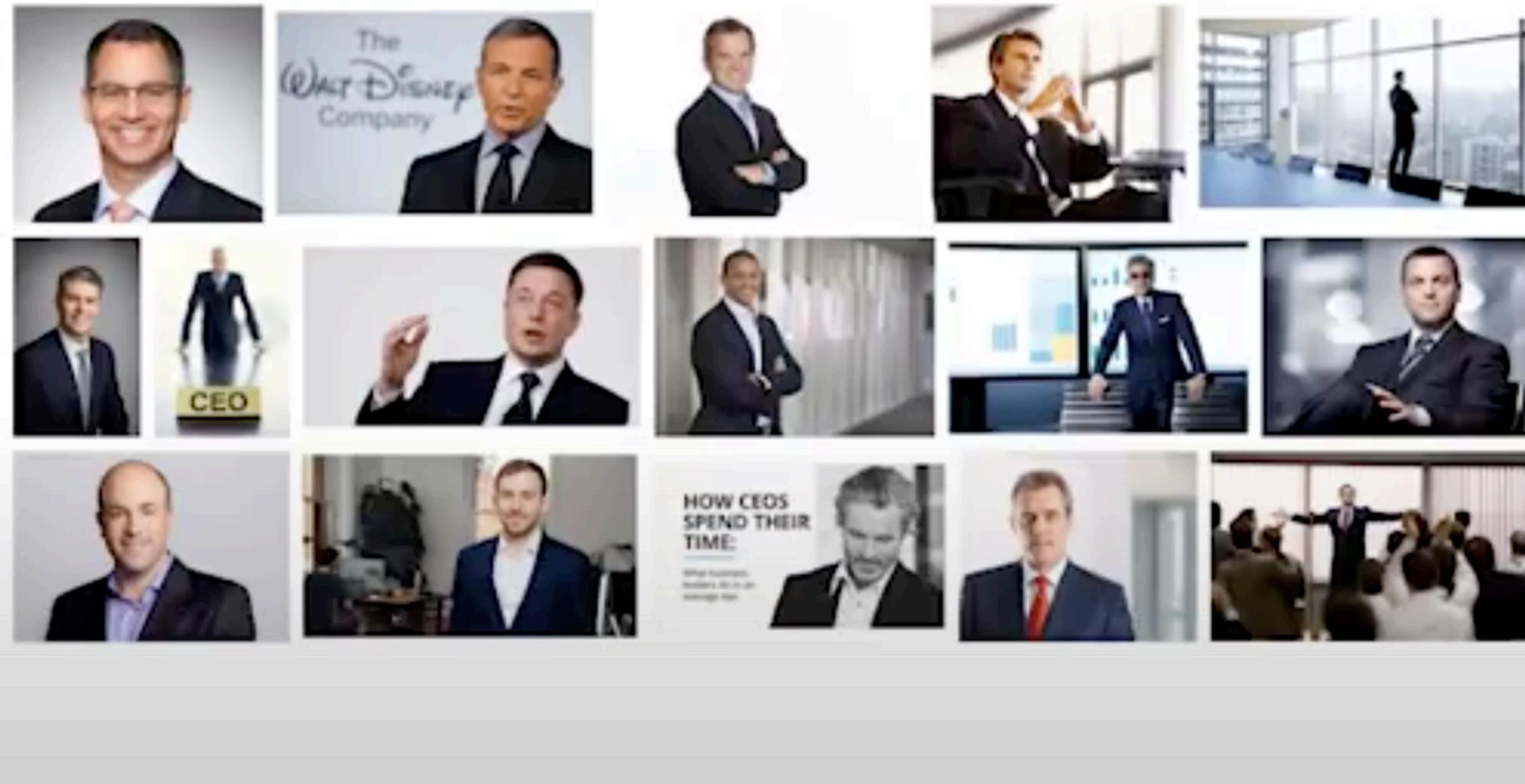
Timnit Gebru, then a graduate student at Stanford University, remembers counting only six Black people other than herself, all of whom she knew, all of whom were men.

“I’m not worried about machines taking over the world. I’m worried about groupthink, insularity and arrogance in the A.I. community — especially with the current hype and demand for people in the field,” she wrote. “The people creating the technology are a big part of the system. If many are actively excluded from its creation, this technology will benefit a few while harming a great many.”

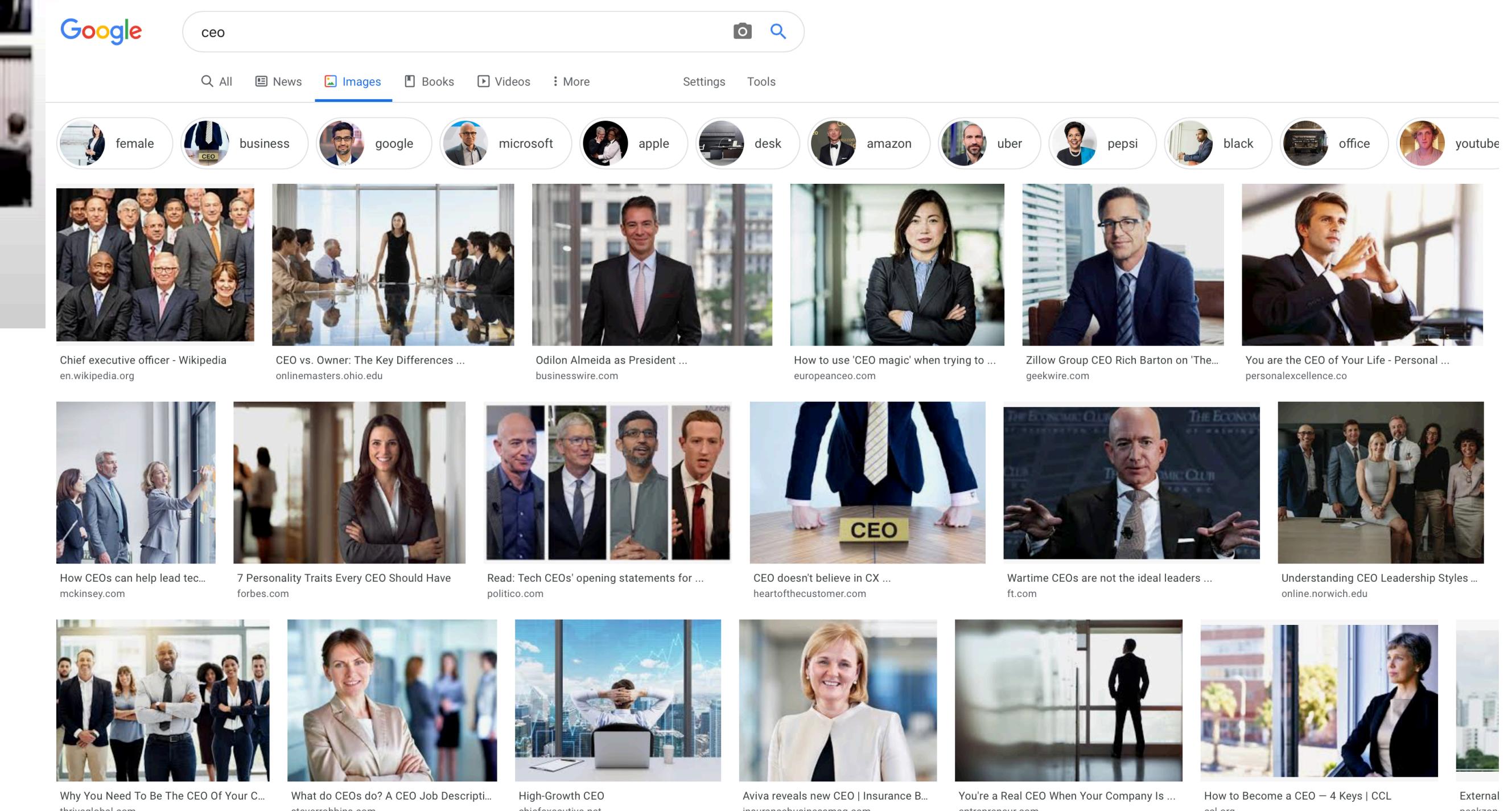
<https://www.nytimes.com/2021/03/15/technology/artificial-intelligence-google-bias.html>

Another example

Image search: CEO



Recent improvement!



Translation

Google translate: English → Turkish → English

The top screenshot shows the input "She is a doctor." and "He is a nurse." being translated to "O bir doktor." and "O bir hemşire." respectively. The bottom screenshot shows the input "O bir doktor." and "O bir hemşire" being translated to "He is a doctor." and "She is a nurse." This illustrates a known issue where Google Translate sometimes fails to correctly handle gendered nouns in one direction.

Recent improvement!

This screenshot shows a more recent version of Google Translate where the same inputs ("O bir doktor." and "O bir hemşire") are correctly translated to "She is a doctor." and "He is a doctor." respectively, indicating a recent improvement in gender handling.

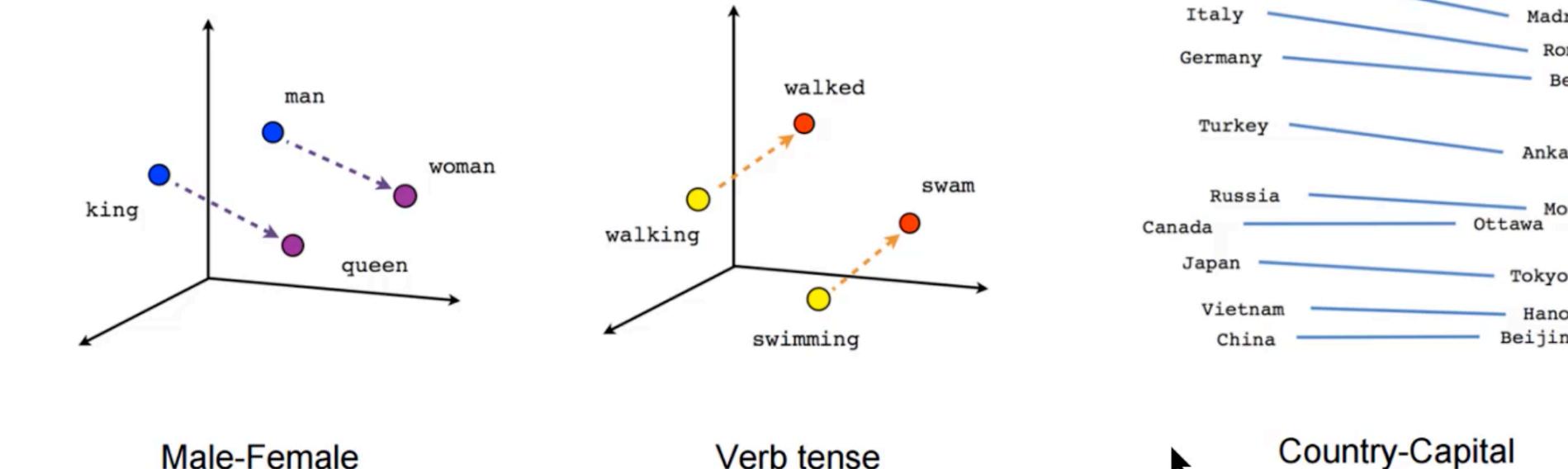
Aravind Narayanan - 21 Fairness Definitions

Word Embeddings



- Word2Vec introduced vector math on word embeddings
- Reveal harmful biases encoded in our language corpora
- Potential solution: de-bias at training time, but at least make user aware

Word analogies are useful



These analogies demonstrate that the embeddings have captured semantic meaning.

Quantifying and Reducing Stereotypes in Word Embeddings

Tolga Bolukbasi¹
Kai-Wei Chang²
James Zou²
Venkatesh Saligrama¹
Adam Kalai²

TOLGAB@BU.EDU
KW@KWCHANG.NET
JAMESZYU@GMAIL.COM
SRV@BU.EDU
ADAM.KALAI@MICROSOFT.COM

MIT Technology Review

Intelligent Machines

How Vector Space Mathematics Reveals the Hidden Sexism in Language

As neural networks tease apart the structure of language, they are finding a hidden gender bias that nobody knew was there.

by Emerging Technology from the arXiv July 27, 2016

Topics+ Top Stories

• Father → Doctor ::
Mother → Nurse

• Man → Computer Programmer ::
Woman → Homemaker

GPT-3

Jerome Pesenti
@an_open_mind

#gpt3 is surprising and creative but it's also unsafe due to harmful biases. Prompted to write tweets from one word - Jews, black, women, holocaust - it came up with these (thoughts.sushant-kumar.com). We need more progress on #ResponsibleAI before putting NLG models in production.

"Jews love money, at least most of the time." "Jews don't read Mein Kampf; they write it."

"#blacklivesmatter is a harmful campaign." "Black is to white as down is to up."

"Women have such a tough time being women. They have periods, do the lifting, and always have to ask for directions." "The best female startup founders are named... Girl."

"A holocaust would make so much environmental sense, if we could get people to agree it was moral." "Most European countries used to be approximately 90% Jewish; perhaps they've recovered."

6:57 AM · Jul 18, 2020 · Twitter Web App

Recent Language Models

On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?

Emily M. Bender*

ebender@uw.edu

University of Washington

Seattle, WA, USA

Timnit Gebru*

timnit@blackinai.org

Black in AI

Palo Alto, CA, USA

Angelina McMillan-Major

aymm@uw.edu

University of Washington

Seattle, WA, USA

Shmargaret Shmitchell

shmargaret.shmitchell@gmail.com

The Aether

- Paper apparently led to authors being fired by Google

We provide recommendations including weighing the environmental and financial costs first, investing resources into curating and carefully documenting datasets rather than ingesting everything on the web, carrying out pre-development exercises evaluating how the planned approach fits into research and development goals and supports stakeholder values, and encouraging research directions beyond ever larger language models.

http://faculty.washington.edu/ebender/papers/Stochastic_Parrots.pdf

Seeing the water

- Whether language models should reflect the world as it *is in the data*, or as society believes it *should be*, depends on what use they are applied to.
- And when we do want them to reflect the world as it *should be*, do we agree on what that is?

Face Recognition

The New York Times

- Old context: no expectation of privacy in public
- New context: "in public" now means "on any street in any city, or on any website on the internet"
- Is it ethical to work on this?
- Is it a problem if it does not work as well on some ethnicities?

The Secretive Company That Might End Privacy as We Know It

A little-known start-up helps law enforcement match photos of unknown people to their online images — and “might lead to a dystopian future or something,” a backer says.



At the same time, debating the merits of these technologies on the basis of their likely accuracy for different groups may distract from a more fundamental question: should we ever deploy such systems, even if they perform equally well for everyone? We may want to regulate the police's access to such tools, even if the tools are perfectly accurate. Our civil rights—freedom of movement and association—are equally threatened by these technologies when they fail and when they work well.

<https://fairmlbook.org/introduction.html>

Questions?

Best practices

What do practitioners need?

2019

- support in fairness-aware data collection and curation
- overcoming teams' blind spots
- implementing more proactive fairness auditing processes
- auditing complex ML systems
- deciding how to address particular instances of unfairness
- addressing biases in the humans embedded throughout the ML development pipeline

Improving Fairness in Machine Learning Systems: What Do Industry Practitioners Need?

Kenneth Holstein
Carnegie Mellon University
Pittsburgh, PA
kjholste@cs.cmu.edu

Jennifer Wortman Vaughan
Microsoft Research
New York, NY
jenn@microsoft.com

Hal Daumé III
Microsoft Research &
University of Maryland
New York, NY
me@hal3.name

Miroslav Dudík
Microsoft Research
New York, NY
mdudik@microsoft.com

Hanna Wallach
Microsoft Research
New York, NY
wallach@microsoft.com

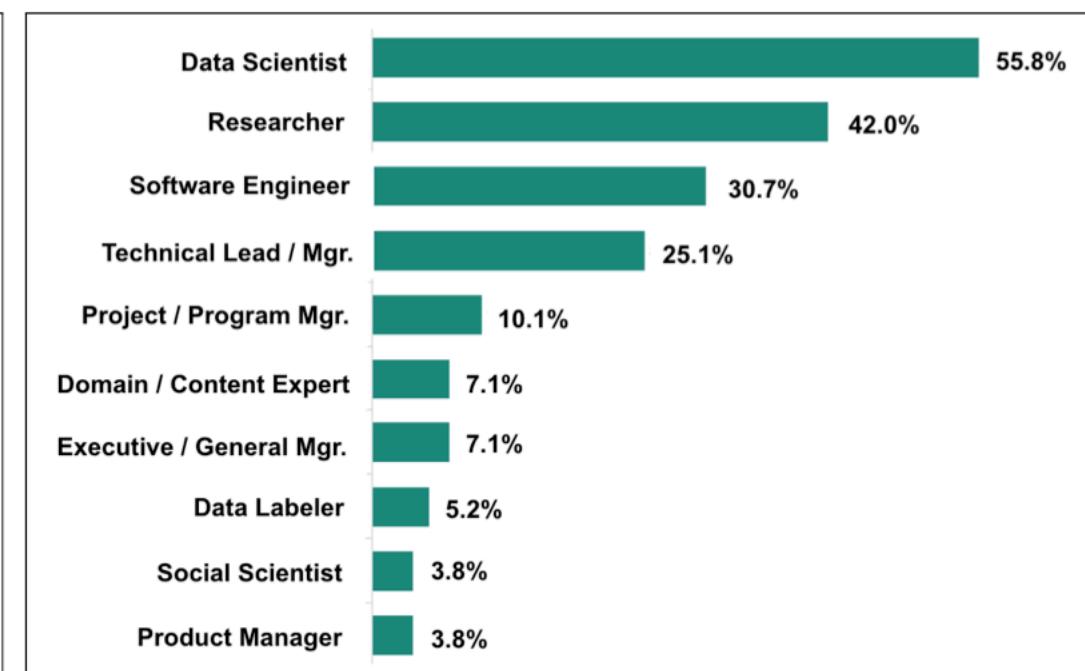
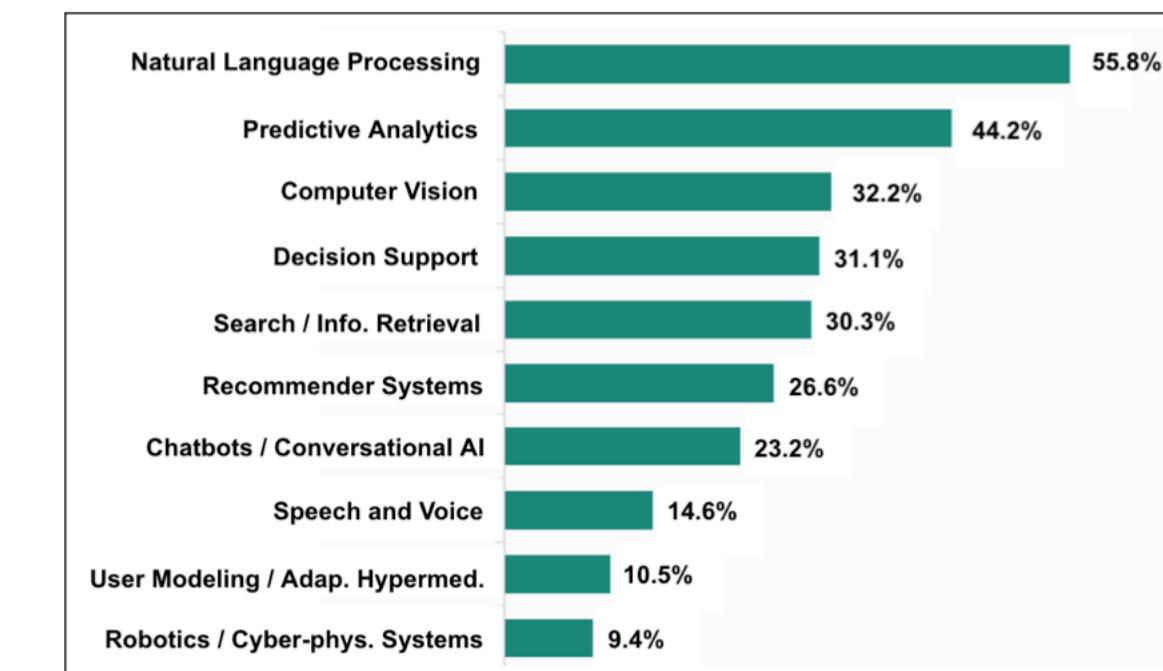


Figure 1: Survey demographics: the top 10 self-reported technology areas (left) and team roles (right).



Rachel Thomas
15.6K subscribers

Some suggestions

- **Ethical risk sweeping**
treat like cybersecurity penetration testing
- **Expanding the ethical circle**
whose interests, desires, experiences, values have we just assumed instead of consulted?
- **Think about the terrible people**
Who might abuse, steal, weaponize what we build? What incentives are we creating?
- **Closing the loop**
Remember that this is not a process to complete and forget. Set up ways to keep improving.

<https://www.youtube.com/watch?v=av7utkFXbU4>

Model Cards



Face Detection

Model Card v0 Cloud Vision API

Overview

Limitations

Trade-offs

Performance

Test your own images

Provide feedback

Explore

Object Detection

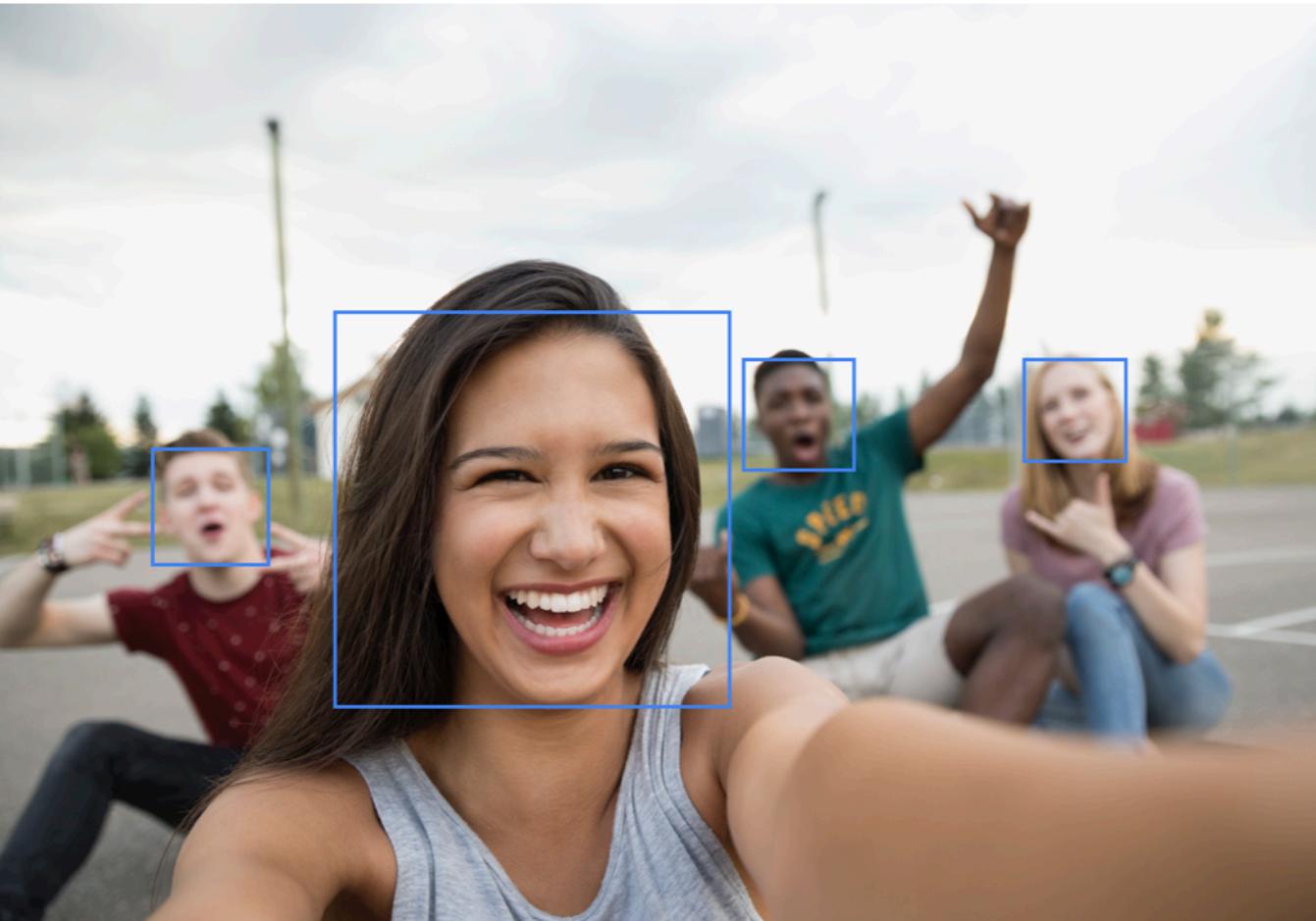
About Model Cards

Face Detection

The [model](#) analyzed in this card detects one or more faces within an image or a video frame, and returns a box around each face along with the location of the faces' major landmarks. The model's goal is exclusively to identify the existence and location of faces in an image. It does not attempt to discover identities or demographics.

On this page, you can learn more about how well the model performs on images with different characteristics, including face demographics, and what kinds of images you should expect the model to perform well or poorly on.

MODEL DESCRIPTION



Input: Photo(s) or video(s)

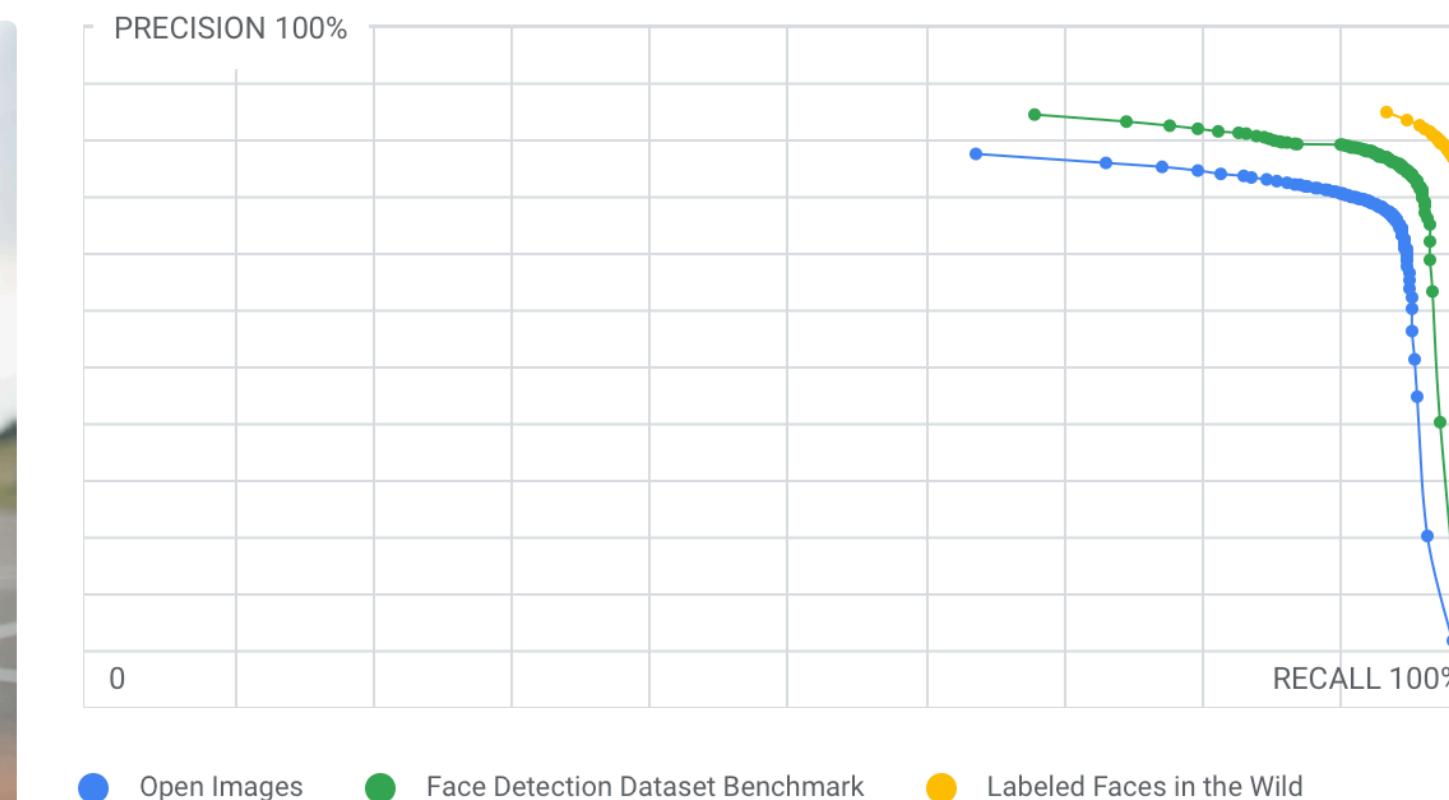
Output: For each face detected in a photo or video, the model outputs:

- [Bounding box](#) coordinates
- Facial landmarks (up to 34 per face)
- Facial orientation (roll, pan, and tilt angles)
- Detection and landmarking confidence scores.

No identity or demographic information is detected.

Model architecture: MobileNet CNN fine-tuned for face detection with a [single shot multibox detector](#).

PERFORMANCE



Overall model performance, and performance [sliced by](#) different image and face characteristics, were assessed, including:

- Derived characteristics (face size, facial orientation, and occlusion)
- Face demographics (human-perceived gender presentation, age, and skin tone)

Overall performance measured with [Precision-Recall \(PR\) values](#) and [Area Under the PR Curve \(PR-AUC\)](#) - standard metrics for evaluating computer vision classifiers. Download raw performance results data [here](#).

Disaggregated performance measured with [Recall](#), which captures how often the model misses faces with specific characteristics. Equal recall across subgroups corresponds to the "[Equality of Opportunity](#)" fairness criterion.

Performance evaluated on: Three research benchmarks distinct from the training set:

<https://modelcards.withgoogle.com/face-detection>

Bias Audit

Center for Data Science and Public Policy



About

Bias Audit Tool

Code

Documentation

Paper



Why we created Aequitas

Machine Learning, AI and Data Science based predictive tools are being increasingly used in problems that can have a drastic impact on people's lives in policy areas such as criminal justice, education, public health, workforce development and social services. Recent work has raised concerns on the risk of unintended bias in these models, affecting individuals from certain groups unfairly. While a lot of bias metrics and fairness definitions have been proposed, there is no consensus on which definitions and metrics should be used in practice to evaluate and audit these systems. Further, there has been very little empirical work done on using and evaluating these measures on real-world problems, especially in public policy.

Aequitas, an open source bias audit toolkit developed by the [Center for Data Science and Public Policy](#) at University of Chicago, can be used to audit the predictions of machine learning based risk assessment tools to understand different types of biases, and make informed decisions about developing and deploying such systems.

Center for Data Science and Public Policy



Bias and Fairness Audit Report

Generated by Aequitas for [Large US City] Criminal Justice Project
January 29, 2018

Project Goal: Identify individuals likely to get booked/charged by police in the near future

Performance Metric: Accuracy (Precision) in the top 150 identified individuals

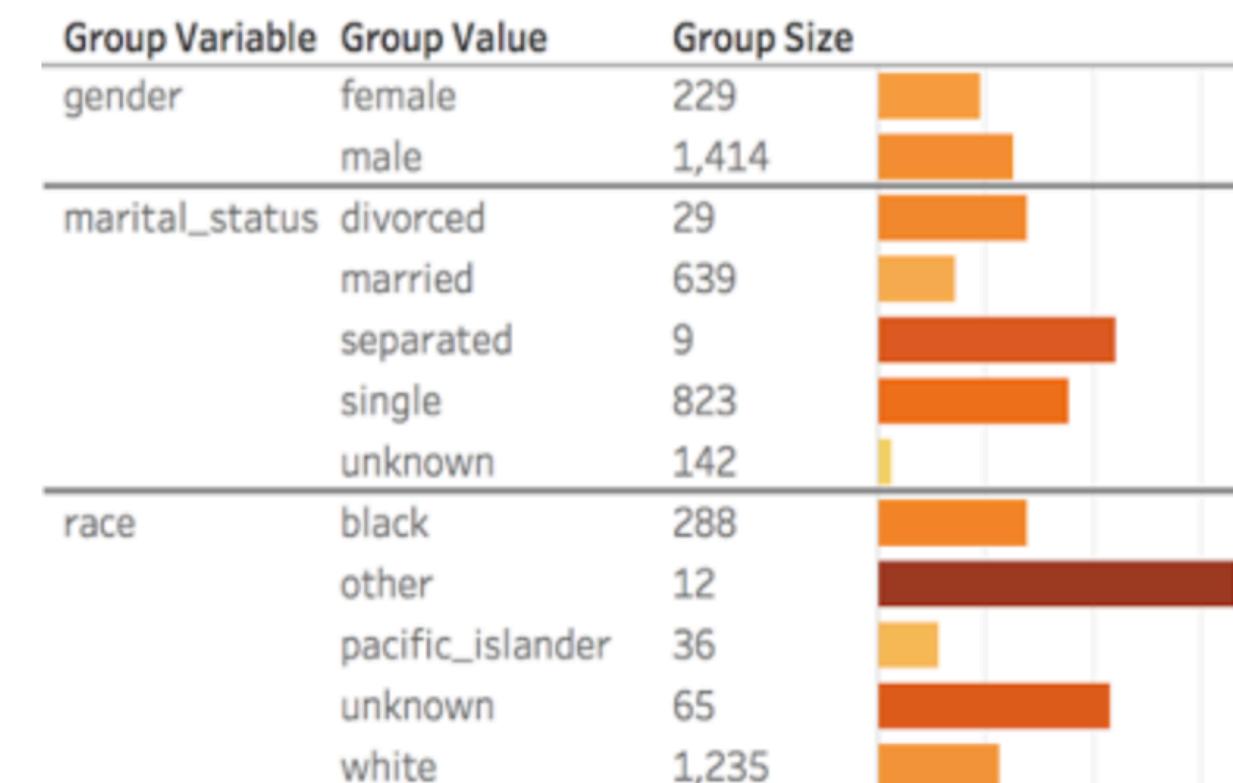
Bias Metrics Considered: Demographic Disparity, Impact Disparity, FPR Disparity, FNR Disparity, FOR Disparity, FDR Disparity

Reference Groups: Race/Ethnicity – White, Gender: Male, Age: None

Model Audited: #841 (Random Forest)

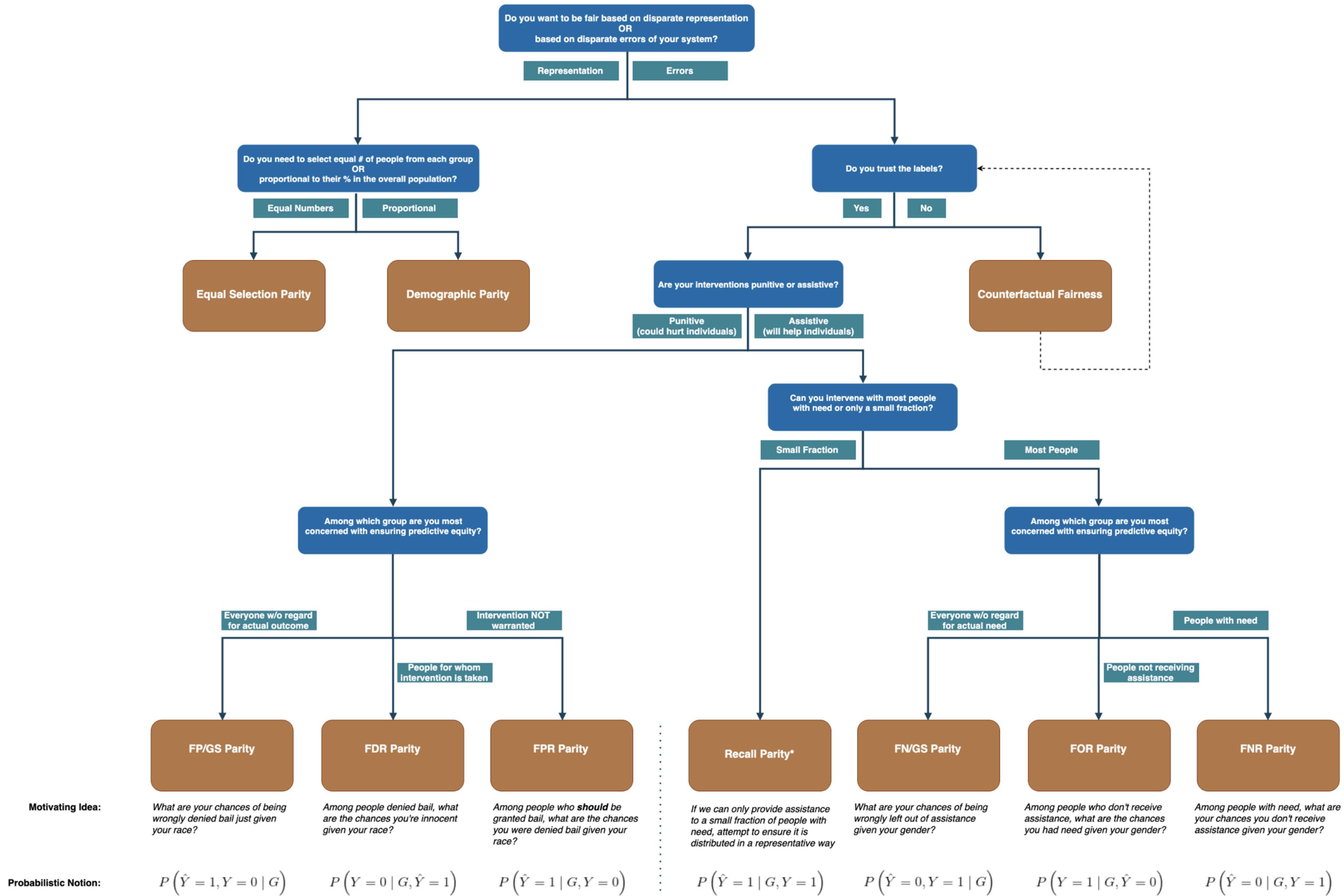
Model Performance: 73%

⚠️ Aequitas has found that Model 841 is **BIASED**. The Bias is in the following attributes:

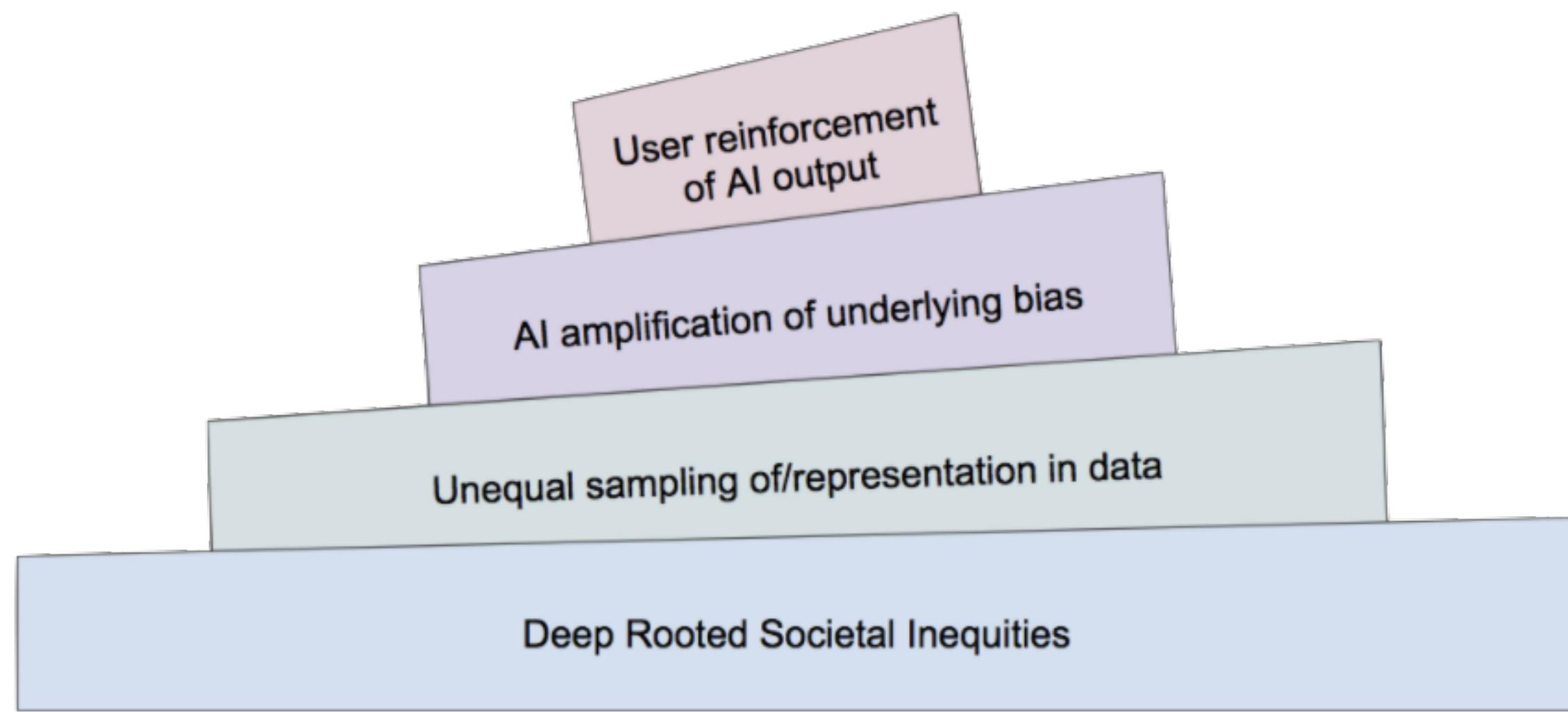


<http://www.datasciencepublicpolicy.org/projects/aequitas/>

FAIRNESS TREE



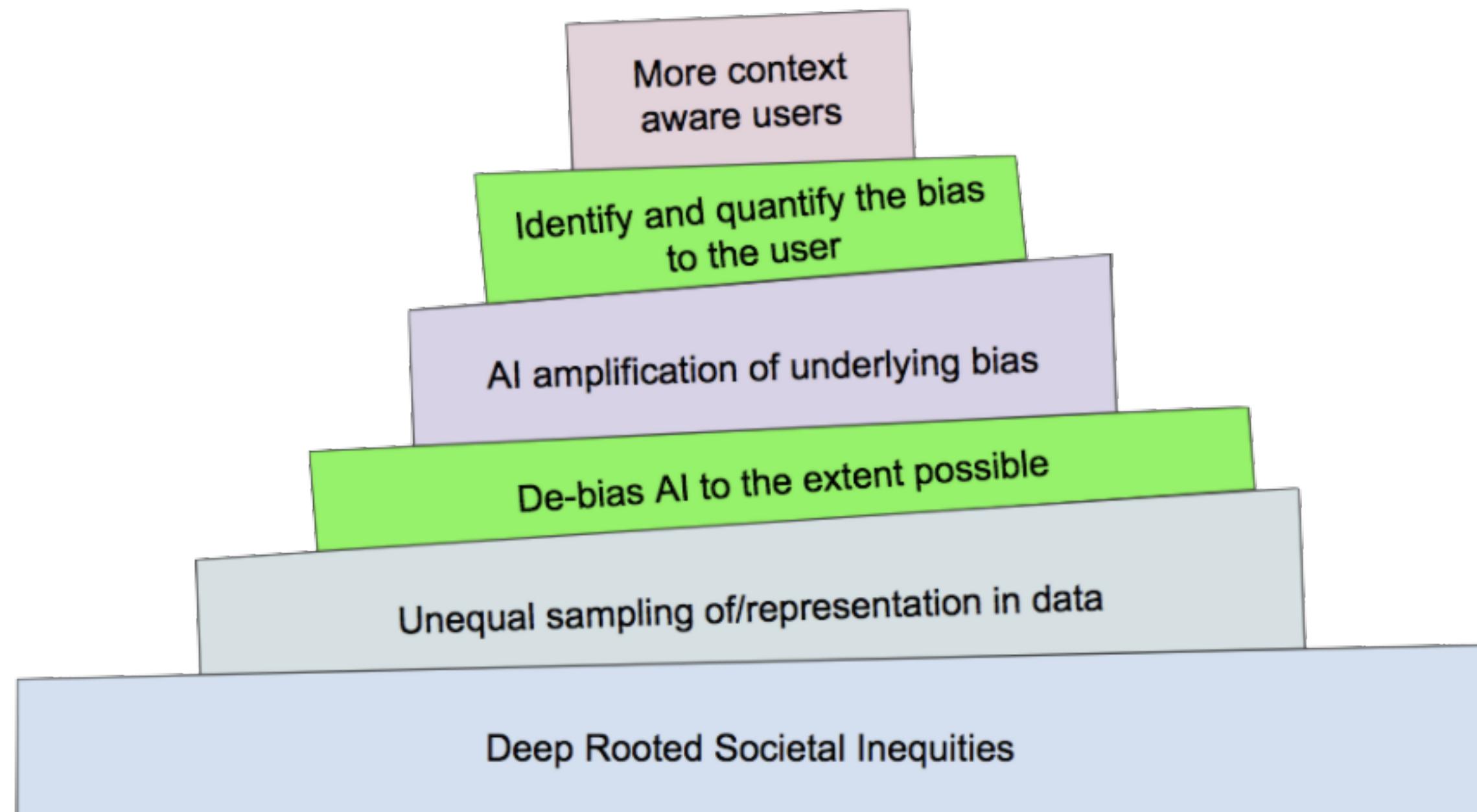
<http://www.datasciencepublicpolicy.org/projects/aequitas/>



AI bias is rooted in societal inequities, manifested in data, and amplified by AI.

Unbiased societal AI is an impossible goal for a single project to achieve, but greatly mitigating its harmful effects is not.

Visualization courtesy of Eric Wang



Display insights about bias contextually and intuitively in product

Make AI fairness and transparency a key component of our work

Build a diverse AI team to bring lived experience to their work

Visualization courtesy of Eric Wang

A professional code of ethics?

Medicine

I will respect the hard-won scientific gains of those physicians in whose steps I walk, and gladly share such knowledge as is mine with those who are to follow.

I will apply, for the benefit of the sick, all measures [that] are required, avoiding those twin traps of overtreatment and **therapeutic nihilism**.

I will remember that there is art to medicine as well as science, and that warmth, sympathy, and understanding may outweigh the surgeon's knife or the chemist's drug.

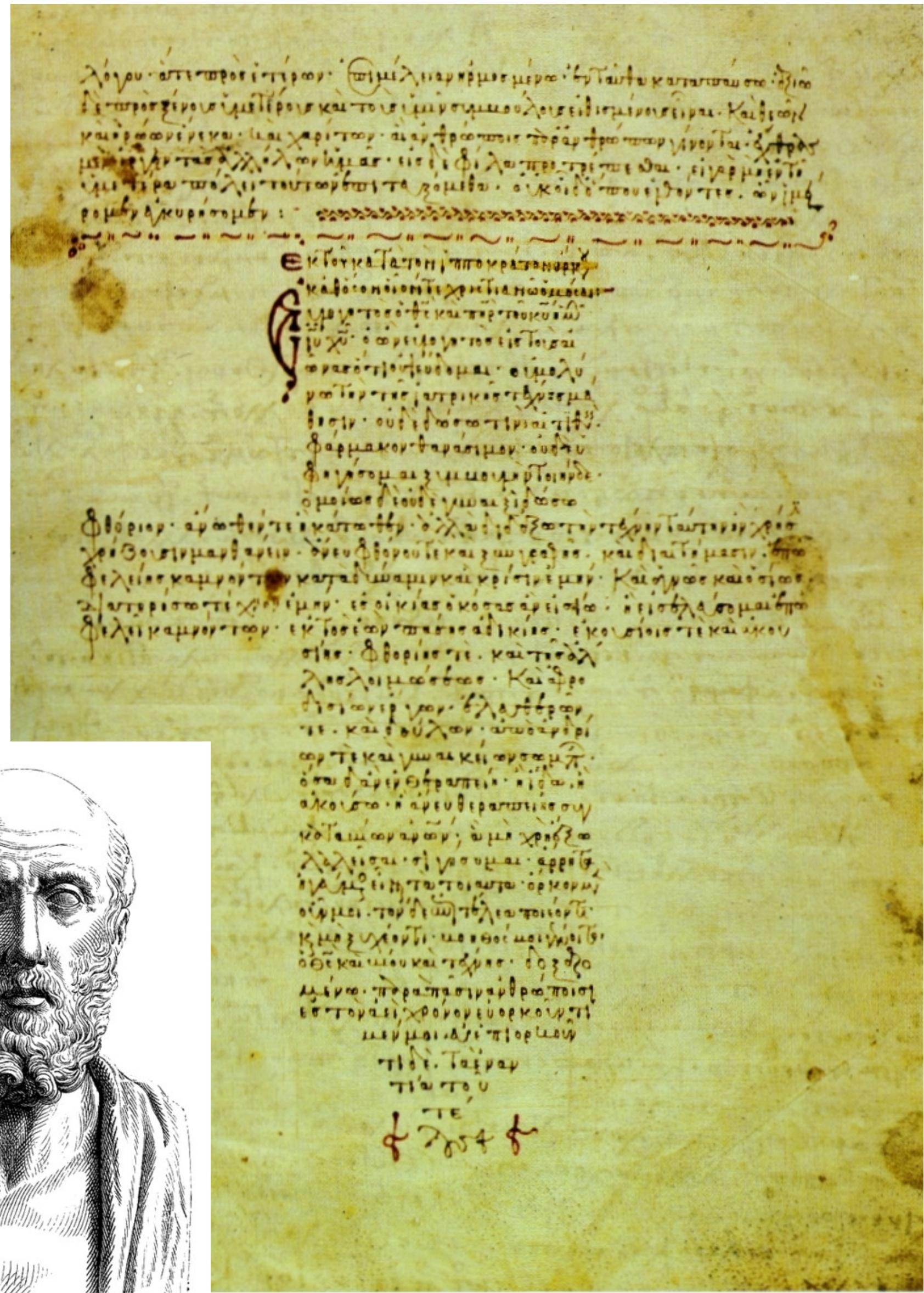
I will not be ashamed to say "I know not", nor will I fail to call in my colleagues when the skills of another are needed for a patient's recovery.

I will respect the privacy of my patients, for their problems are not disclosed to me that the world may know. Most especially must I tread with care in matters of life and death. If it is given me to save a life, all thanks. But it may also be within my power to take a life; this awesome responsibility must be faced with great humbleness and awareness of my own frailty. Above all, I must not play at God.

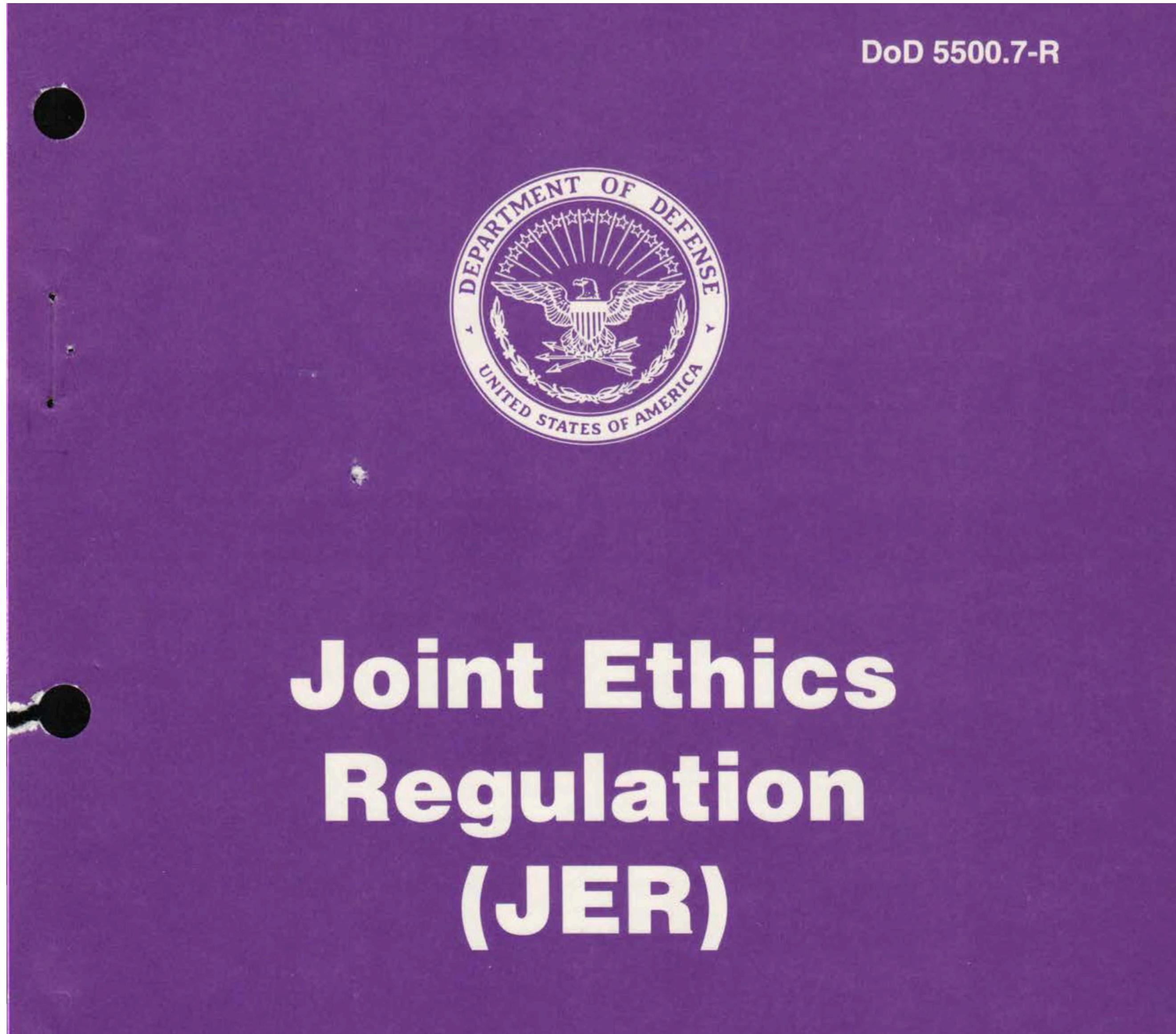
I will remember that I do not treat a fever chart, a cancerous growth, but a sick human being, whose illness may affect the person's family and economic stability. My responsibility includes these related problems, if I am to care adequately for the sick.

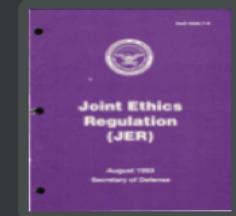
I will prevent disease whenever I can, for prevention is preferable to cure.

I will remember that I remain a member of society, with special obligations to all my fellow human beings, those sound of mind and body as well as the infirm.



Armed Forces





1



2



3



4



5



6



7



8



9



10



11



12



13



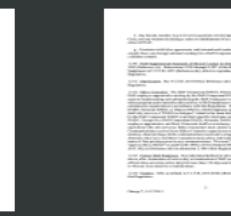
14



15



16



17



18



19



20



21



22



23



24



25



26



27



28



29



30



31



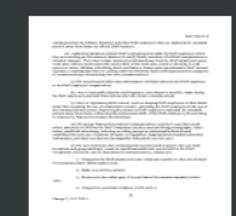
32



33



34



35



36



37



38



39



40



41



42



43



44



45



46



47



48



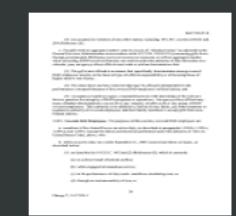
49



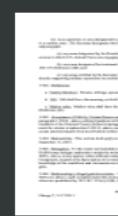
50



51



52



53



54



55



56



57



58



59



60



61



62



63



64



65



66



67



68



69



70



71



72



73



74



75



76



77



78



79



80



81



82



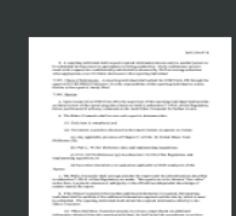
83



84



85



86



87



88



89



90



91



92



93



94



95



96



97



98



99



Questions?

Learn more

Experts

- Practical Data Ethics course from Rachel Thomas: <https://ethics.fast.ai>
- Single lecture from Fast.ai 2020 course: <https://www.youtube.com/watch?v=krIVOb23EH8>



Rachel Thomas
@math_rachel
Director of USF Center for Applied Data Ethics @DataInstituteSF + co-founder fast.ai | deep learning, ethics, math PhD | photo by Gabriela Hasbun | she/her
📍 San Francisco, CA ⚡ fast.ai/topics/#ai-in-... 📅 Joined May 2013
694 Following 78.7K Followers

Topics covered:

1. Disinformation
2. Bias & Fairness
3. Ethical Foundations & Practical Tools
4. Privacy & surveillance
5. Our Ecosystem: Metrics, Venture Capital, & Losing the Forest for the Trees
6. Algorithmic Colonialism, and Next Steps

Fairness and machine learning

Limitations and Opportunities

Solon Barocas, Moritz Hardt, Arvind Narayanan

This online textbook is an incomplete work in progress. Essential chapters are still missing. In the spirit of open review, we solicit broad feedback that will influence existing chapters, as well as the development of later material.

CONTENTS

ABOUT THIS BOOK

1 INTRODUCTION

[PDF](#)

2 CLASSIFICATION

[PDF](#)

We introduce formal non-discrimination criteria, establish their relationships, and illustrate their limitations.

3 LEGAL BACKGROUND AND NORMATIVE QUESTIONS

We survey the literature on discrimination in law, sociology, and philosophy. We then discuss the challenges that arise in translating these ideas of fairness to the statistical decision-making setting.

4 CAUSALITY

[PDF](#)

We dive into the rich technical repertoire of causal inference and how it helps articulate and address shortcomings of the classification paradigm, while raising new conceptual and normative questions.

5 TESTING DISCRIMINATION IN PRACTICE

[PDF](#)

We systematize tests of discrimination and discuss the practical complexities of applying them, both to traditional decision-making systems and to algorithmic systems.

Note: For an updated resource, please see fairmlbook.org.

CS 294: Fairness in Machine Learning

UC Berkeley, Fall 2017

Time: Monday and Friday 2:30PM - 3:59PM

Location: Soda 405

Instructor: Moritz Hardt

<https://fairmlbook.org>

Dealing with Bias and Fairness in Building Data Science/ML/AI Systems

A Hands-on Tutorial

[View on GitHub](#)

Dealing with Bias and Fairness in Data Science Systems: A Practical Hands-on Tutorial

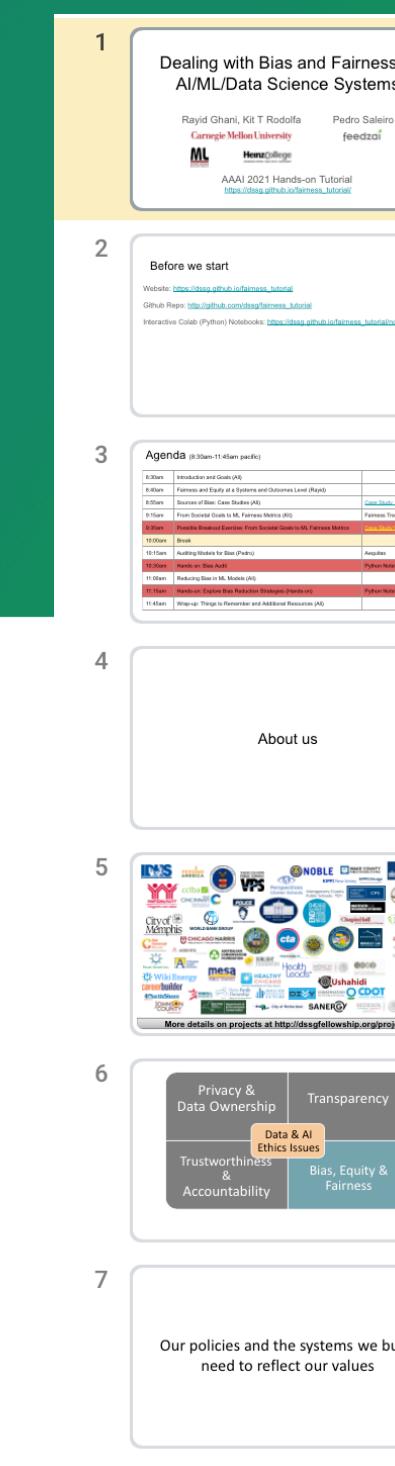
Earlier versions:

- Presented at KDD 2020: [Github repo](#), [Video](#), and [web page](#)

Presenters

- Pedro Saleiro, Feedzai
- Kit T. Rodolfa, Carnegie Mellon University
- Rayid Ghani, Carnegie Mellon University

https://dssg.github.io/fairness_tutorial/



Dealing with Bias and Fairness in AI/ML/Data Science Systems

Rayid Ghani, Kit T Rodolfa

Carnegie Mellon University

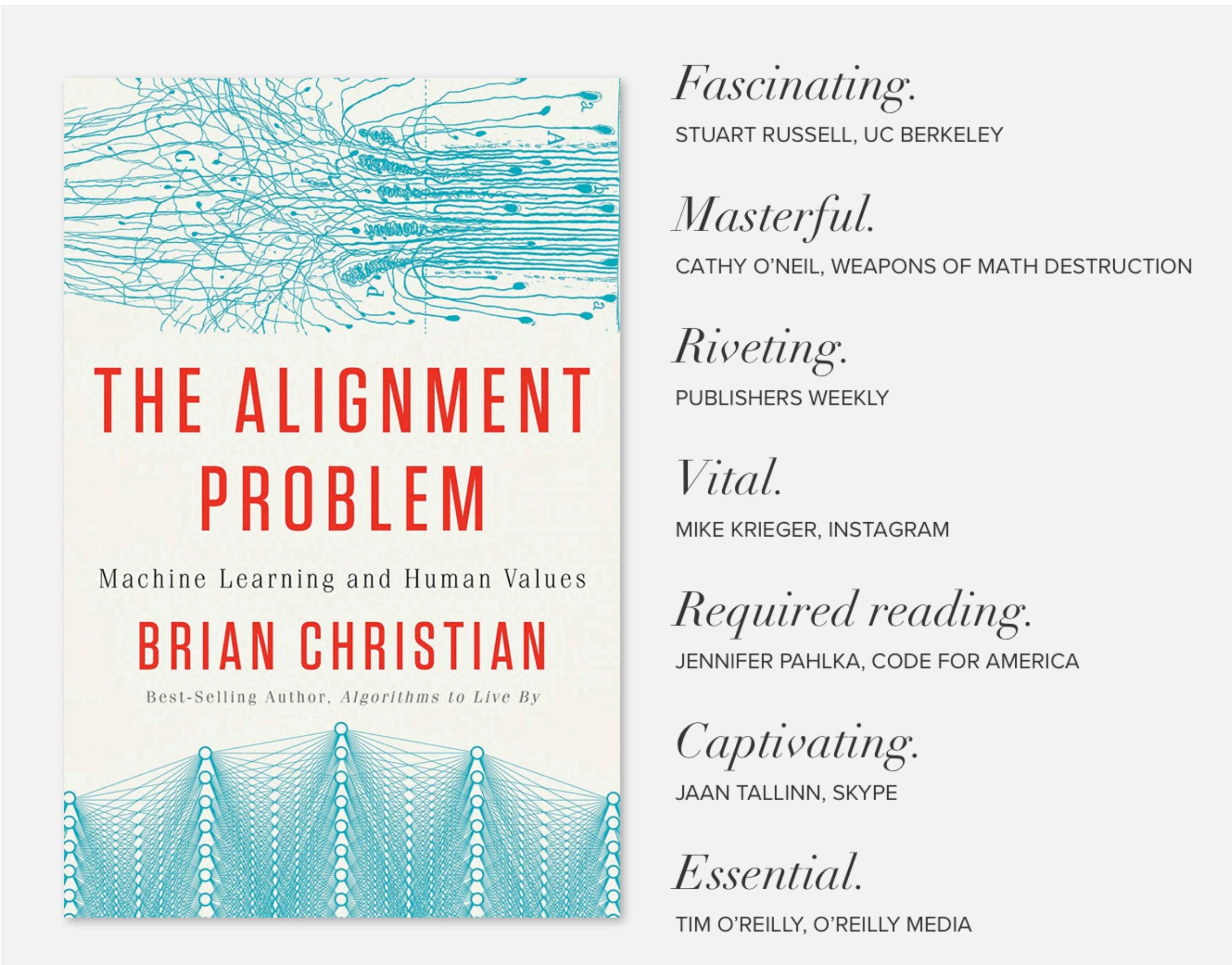


HeinzCollege

INFORMATION SYSTEMS • PUBLIC POLICY • MANAGEMENT

AAAI 2021 Hands-on Tutorial

https://dssg.github.io/fairness_tutorial/



Fascinating.

STUART RUSSELL, UC BERKELEY

Masterful.

CATHY O'NEIL, WEAPONS OF MATH DESTRUCTION

Riveting.

PUBLISHERS WEEKLY

Vital.

MIKE KRIEGER, INSTAGRAM

Required reading.

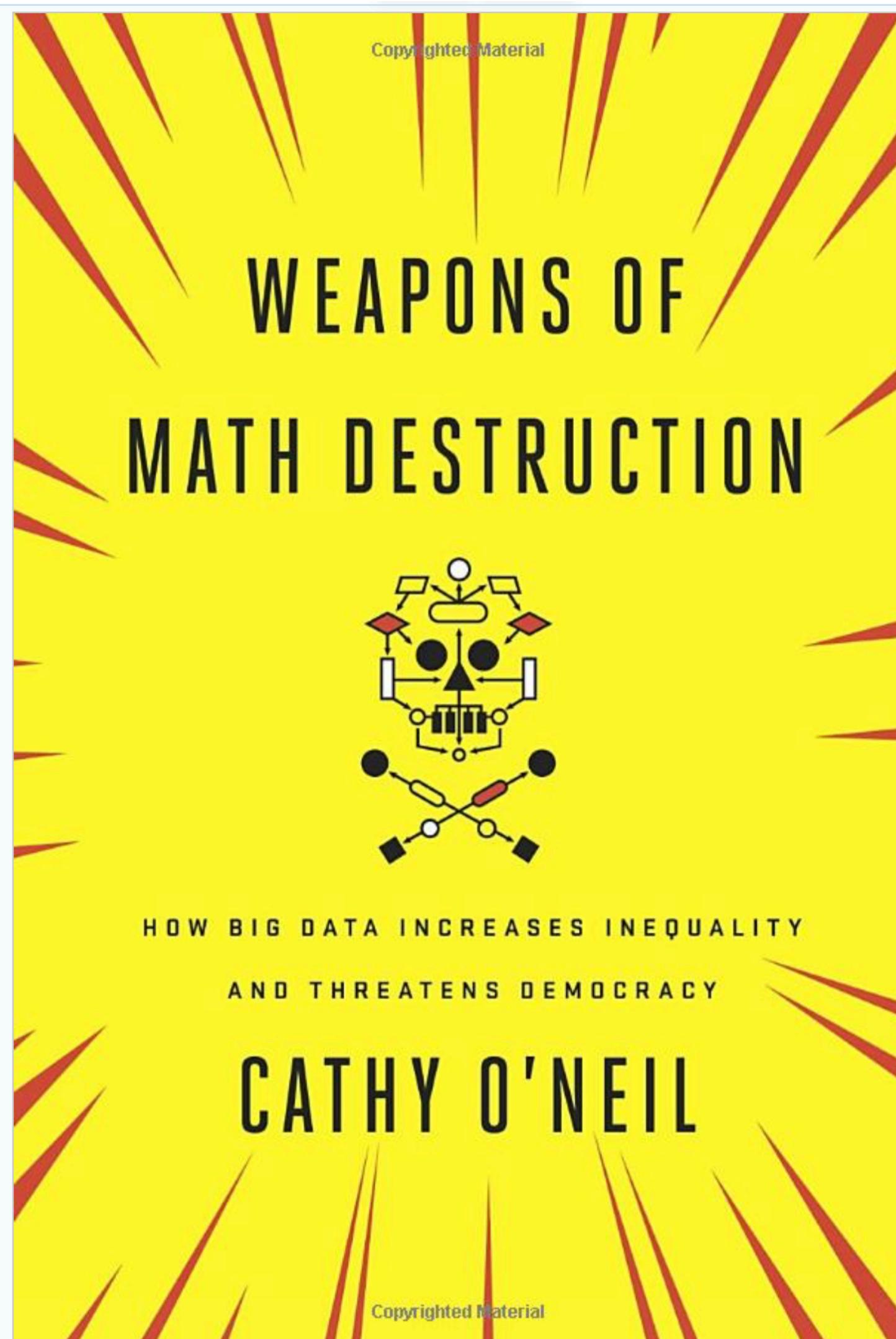
JENNIFER PAHLKA, CODE FOR AMERICA

Captivating.

JAAN TALLINN, SKYPE

Essential.

TIM O'REILLY, O'REILLY MEDIA



Thank you!