# A STORY OF DISCRIMINATIN AND UNFAIRNESS:

# PREJUDICE IN WORD EMBEDDINGS

## Aylin Caliskan    @aylin_cim

Princeton University
CITP Fellow and Postdoctoral Research Associate

# Thanks to:

- Organizers
- Angels
- Chaos mentors (did you know that they existed?)
- Assemblies
- Artists
- CCC
  - Programmer de-anonymization
  - Stylometry

# Thanks to my co-authors!

Joanna Bryson

@j2bryson

Arvind Narayanan

@random_walker

# A new approach to algorithmic transparency

– Not about classification unfairness discovery

– Uncovering societal bias embedded in machine learning models for:

- Machine translation
- Sentiment analysis: market trends - company reviews, customer satisfaction - movie reviews…
- Web search and search engine optimization hacks
  - Filter bubble

# Disclaimer:

Examples with offensive content.
Does not reflect our opinions!

# Problem

- Machine learning models trained on human data.

- Consequently, models reflect human culture and semantics.

- Human culture happens to include:

  – Bias and prejudice

# Problem

- Machine learning models trained on human data.

- Consequently, models reflect human culture and semantics.

- Human culture happens to include:

  - Bias and prejudice $\rightarrow$ unfairness and discrimination ☹

# Problem

- We focus on language models.
- Language models represent semantic spaces with <u>word embeddings</u>

$word_1$,          $feature_1$ , $feature_2$ , $feature_3$ , $feature_4$ , … $feature_{300}$

$word_2$,          $feature_1$ , $feature_2$ , $feature_3$ , $feature_4$ , … $feature_{300}$

$word_3$,          $feature_1$ , $feature_2$ , $feature_3$ , $feature_4$ , … $feature_{300}$

…

$word_{2000000}$, $feature_1$ , $feature_2$ , $feature_3$ , $feature_4$ , … $feature_{300}$

# Problem

- We focus on language models.

- Language models represent semantic spaces with word embeddings

— Meaning

— Syntax

— Similarities

- Woman to man is girl to boy

# Problem

- We focus on language models.
- Language models represent semantic spaces with word embeddings

  – Meaning

  – Syntax

  – Similarities

    - Woman to man is girl to boy

# Problem

- We focus on language models.

- Language models represent semantic spaces with word embeddings

  – Meaning

  – Syntax

  – Similarities

    - Woman to man is girl to boy

    - Paris to France is Rome to Italy

    - Banana to bananas is nut to nuts

# Generating language models



Donald J. Trump ✓
@realDonaldTrump                    ⚙   + Follow

Sadly, because president Obama has done such a poor job as president, you won't see another black president for generations!

RETWEETS   FAVORITES
8,875      7,690

3:15 AM - 25 Nov 2014

# Generating language models

# Generating language models

# Generating language models

# Generating language models

# Generating language models

# Generating language models



word2vec

glove

# Models used in:

- Text generation
- Automated speech generation
- Machine translation
- Sentiment analysis
- Named entity recognition
- Web search…

# Natural language processing as a service:

# Future of AI

**Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours**

# Future of AI



**Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours**

# Future of AI



Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours

# Future of AI



**Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours**

# Future of AI



**Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours**

# Future of AI

# Future of AI



Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours

# Future of AI



Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours

# Future of AI



Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours

# Stereotype threat

Groups: Black and white Americans

Threat: Intellectual ability



**Effects of Stereotype Threat**

"The Effects of Stereotype Threat on the Standardized Test Performance of College Students (adjusted for group differences on SAT)". From J. Aronson, C.M. Steele, M.F. Salinas, M.J. Lustina, *Readings About the Social Animal*, 8th edition, ed. E. Aronson

# Stereotype threat

Groups: Men and women

Threat: Math ability



Stereotype threat and test performance

The effect of stereotype threat (ST) on math test scores for girls and boys. Data from Osborne (2007).[18]

# What to do?

- "Be aware of bias in life. We are constantly being primed.

- Debias by presenting positive alternatives.

- Engage in proactive affirmative efforts not only on the cultural level but also the structural level."

<div style="text-align: right; color: red;">Banaji and Greenwald</div>

# What to do?

- "Be aware of bias in life. We are constantly being primed.

- Debias by presenting positive alternatives.

- Engage in proactive affirmative efforts not only on the cultural level but the structural level."

Algorithmic transparency

Banaji and Greenwald

# What to do?

- "Be aware of bias in life. We are constantly being primed.

- Debias by presenting positive alternatives.

- Engage in proactive affirmative efforts not only on the cultural level but the structural

Algorithmic transparency

Quantify bias in models

eenwald

# How to measure bias?
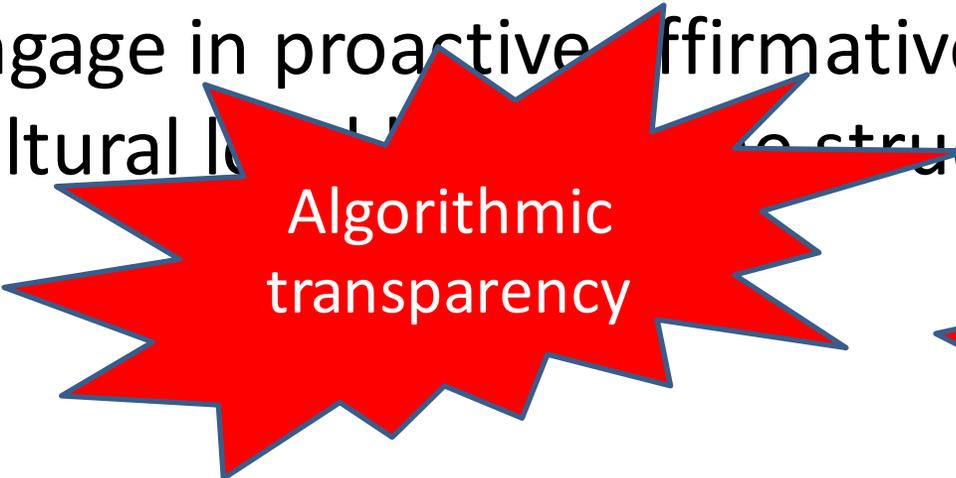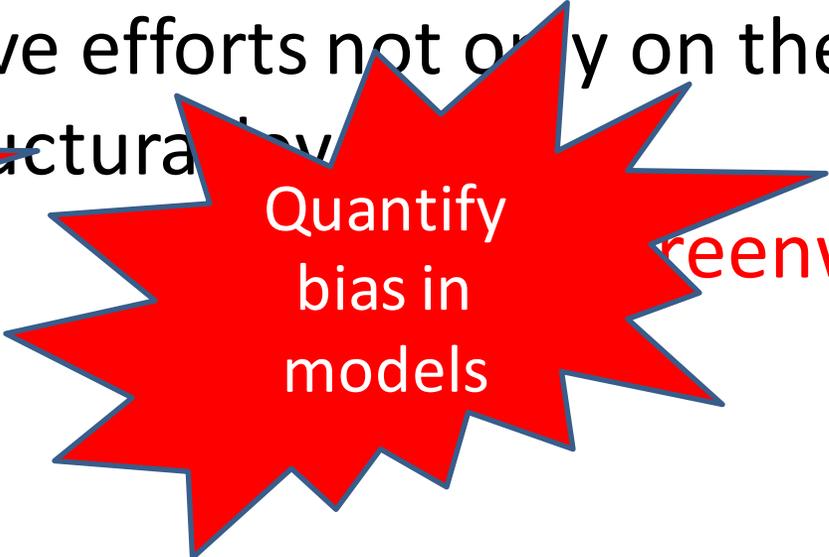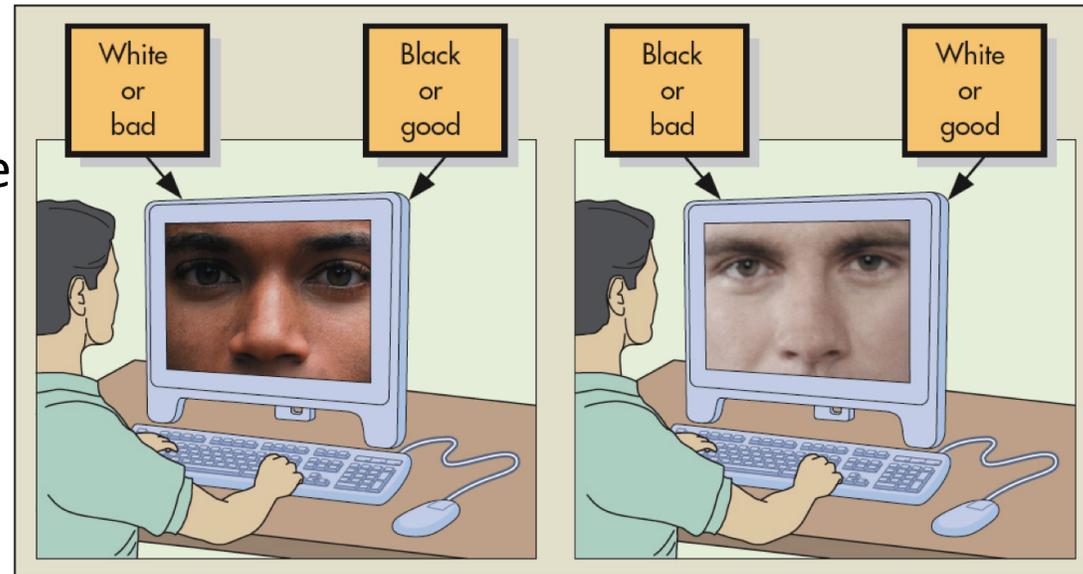
- Implicit Association Test – Greenwald et al. 1998

- Reveals subconscious bias
  - that you might be unaware

- Association of
  - Societal groups
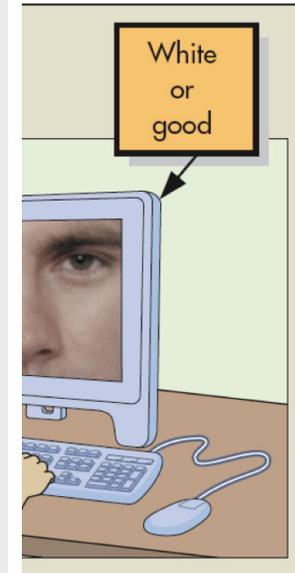    with
  - Stereotype words

# How to measure bias?



## Project Implicit®

LOG IN    TAKE A TEST    ABOUT US    EDUCATION    HELP    CONTACT US    DONATE

- Implicit

- Reveals
  - that

- Associa
  - Socie
      v
  - Stere

| | |
|---|---|
| **Presidents IAT** | *Presidents* ('Presidential Popularity' IAT). This IAT requires the ability to recognize photos of Barack Obama and one or more previous presidents. |
| **Skin-tone IAT** | *Skin-tone* ('Light Skin - Dark Skin' IAT). This IAT requires the ability to recognize light and dark-skinned faces. It often reveals an automatic preference for light-skin relative to dark-skin. |
| **Sexuality IAT** | *Sexuality* ('Gay - Straight' IAT). This IAT requires the ability to distinguish words and symbols representing gay and straight people. It often reveals an automatic preference for straight relative to gay people. |
| **Arab-Muslim IAT** | *Arab-Muslim* ('Arab Muslim - Other People' IAT). This IAT requires the ability to distinguish names that are likely to belong to Arab-Muslims versus people of other nationalities or religions. |
| **Gender-Science IAT** | *Gender - Science.* This IAT often reveals a relative link between liberal arts and females and between science and males. |
| **Native IAT** | *Native American* ('Native - White American' IAT). This IAT requires the ability to recognize White and Native American faces in either classic or modern dress, and the names of places that are either American or Foreign in origin. |
| **Gender-Career IAT** | *Gender - Career.* This IAT often reveals a relative link between family and females and between career and males. |
| **Weight IAT** | *Weight* ('Fat - Thin' IAT). This IAT requires the ability to distinguish faces of people who are obese and people who are thin. It often reveals an automatic preference for thin people relative to fat people. |

**https://implicit.harvard.edu/implicit**

# Measuring bias in Germany



**https://implicit.harvard.edu/germany**

# How do we measure bias?

- **Word Embedding Association Test (WEAT)**
  - Calculate implicit associations between societal categories and evaluative attributes
    - Effect size of bias

# How do we measure bias?

- **Word Embedding Association Test (WEAT)**
  - Calculate implicit associations between societal categories and evaluative attributes
    - Effect size of bias $\dfrac{\text{mean}_{x \in X}\, s(x,A,B) - \text{mean}_{y \in Y}\, s(y,A,B)}{\text{std-dev}_{w \in X \cup Y}\, s(w,A,B)}$

$$s(X,Y,A,B) = \sum_{x \in X} s(x,A,B) - \sum_{y \in Y} s(y,A,B)$$

$$s(w,A,B) = \text{mean}_{a \in A}\cos(\vec{w},\vec{a}) - \text{mean}_{b \in B}\cos(\vec{w},\vec{b})$$

# How do we measure bias?

- **Word Embedding Association Test (WEAT)**
  - Calculate implicit associations between societal categories and evaluative attributes
    - Effect size of bias
    
    $$\frac{\text{mean}_{x \in X}\, s(x, A, B) - \text{mean}_{y \in Y}\, s(y, A, B)}{\text{std-dev}_{w \in X \cup Y}\, s(w, A, B)}$$

    $$s(X, Y, A, B) = \sum_{x \in X} s(x, A, B) - \sum_{y \in Y} s(y, A, B)$$

    $$s(w, A, B) = \text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})$$

    - Statistical significance

    $$\text{Pr}_i[s(X_i, Y_i, A, B) > s(X, Y, A, B)] \quad \textit{where Pr}_i = \textit{null hypothesis}$$

# How do we measure bias?

- **Word Embedding Factual Association Test (WEFAT)**
  - Evaluate association of certain words with specific bias

# How do we measure bias?

- **Word Embedding Factual Association Test (WEFAT)**
  - Evaluate association of certain words with specific bias

$$s(w, A, B) = \frac{\text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})}{\text{std-dev}_{x \in A \cup B} \cos(\vec{w}, \vec{x})}$$

# Baseline: Women with androgynous names



United States™ Census Bureau

## Genealogy

### Frequently Occurring Surnames from Census 1990 – Names Files

Tweet | Share

*NOTE: No specific individual information is given.*

**Files**

TXT   dist.all.last [<1.0MB]
TXT   dist.female.first [<1.0MB]
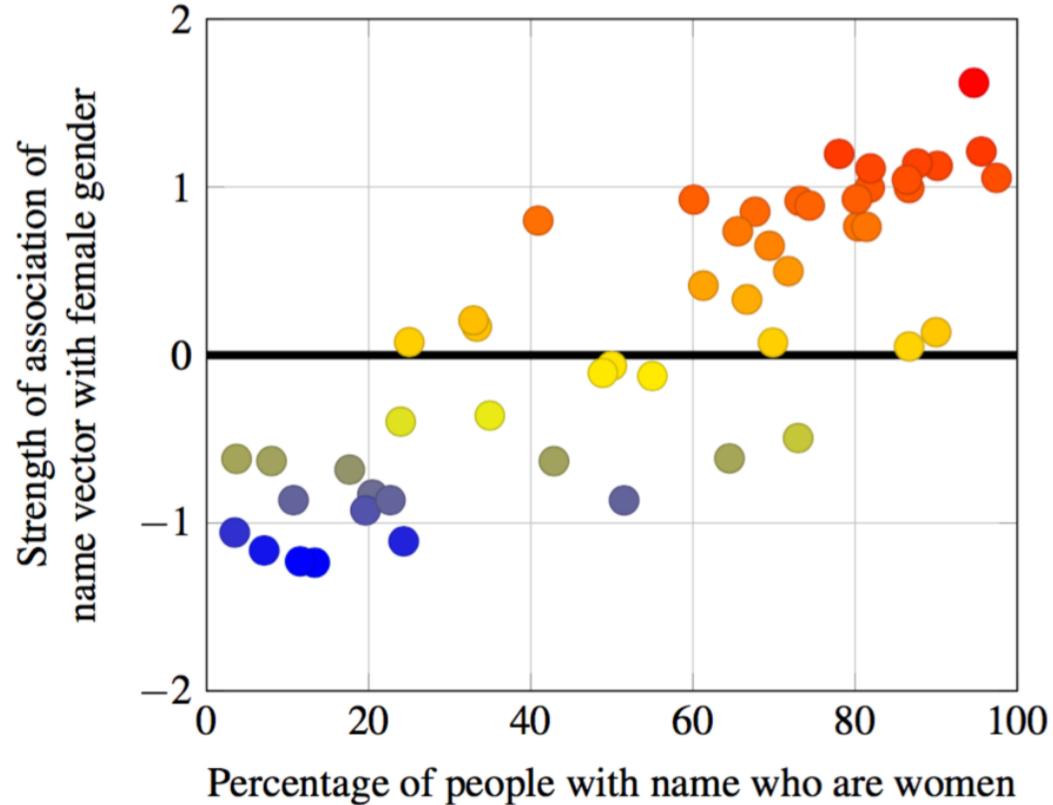TXT   dist.male.first [<1.0MB]

Each of the three files, (dist.all.last), (dist. male.first), and (dist female.first) contain four items of data. The four items are:

1. A "Name"
2. Frequency in percent
3. Cumulative Frequency in percent
4. Rank

In the file (dist.all.last) one entry appears as:
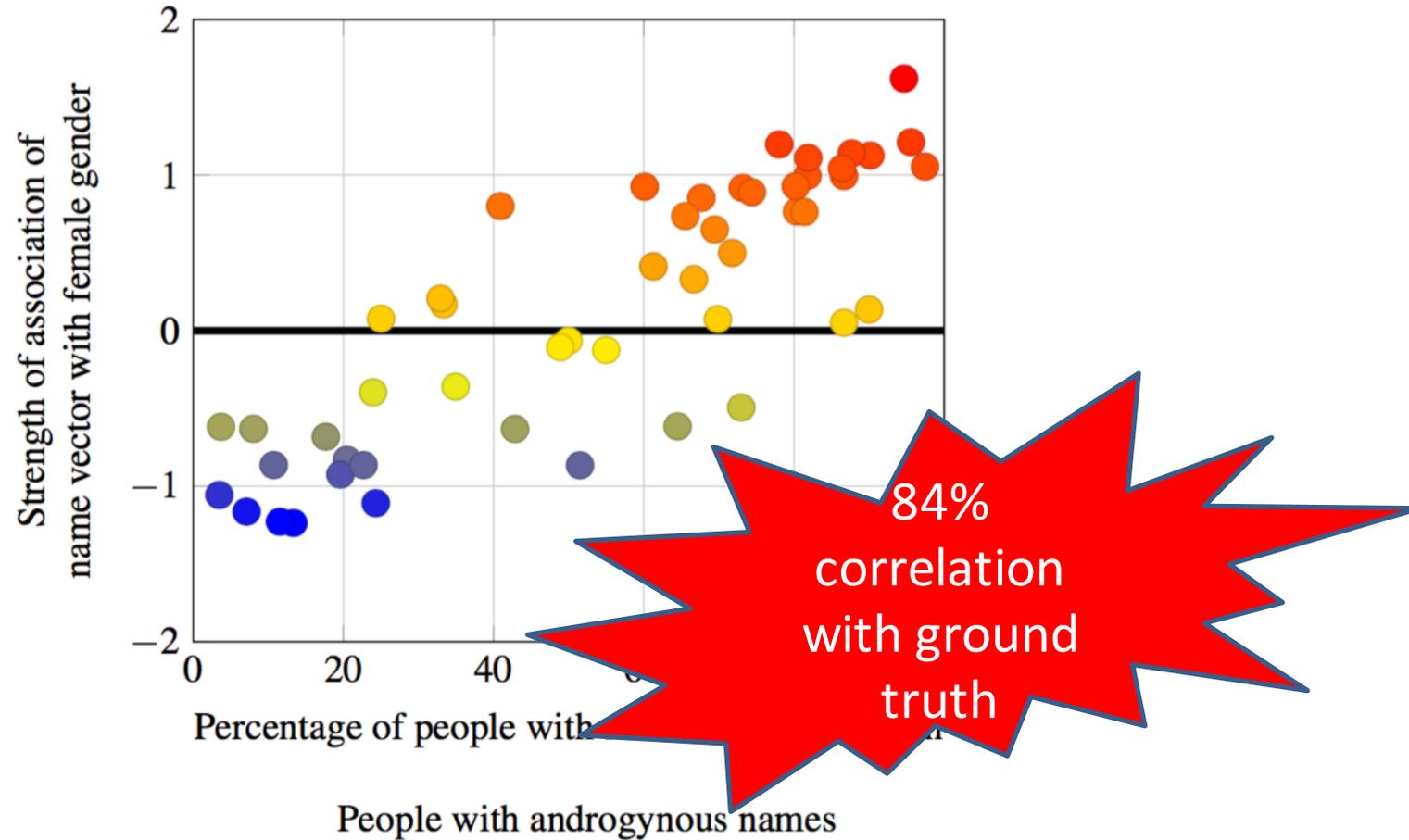
# WEFAT: Women with androgynous names



People with androgynous names
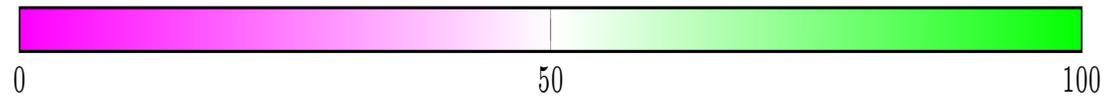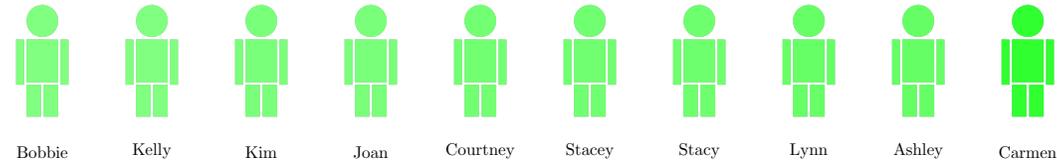Pearson's correlation coefficient $\rho = 0.84$ with $p$-value $< 10^{-13}$.

# WEFAT: Women with androgynous names



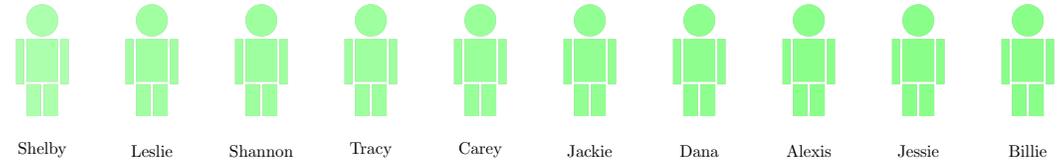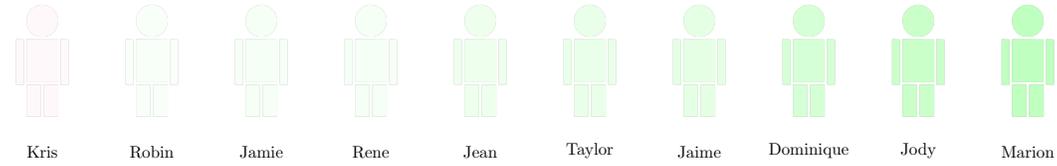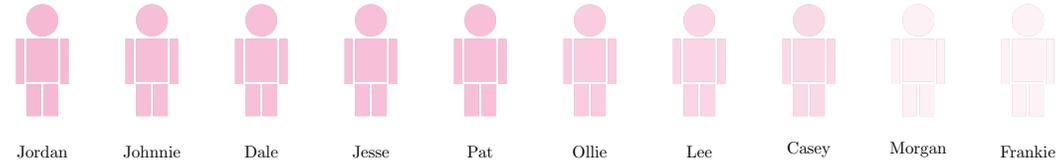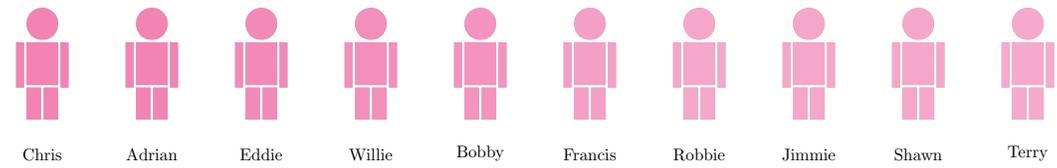People with androgynous names
Pearson's correlation coefficient $\rho = 0.84$ with $p$-value $< 10^{-13}$.

84% correlation with ground truth

Predicted Percentage of Women with Name

Pearson's correlation coefficient $\rho = 0.84$ with 1990 U.S. Census Name and Gender Statistics

# Baseline: Women employed in the US



UNITED STATES DEPARTMENT OF LABOR

A to Z Index | FAQs | About BLS | Contact Us | Subscribe to E-mail Updates | GO

**BUREAU OF LABOR STATISTICS**

Follow Us | What's New | Release Calendar | Blog

Search BLS.gov

Home ▾ | Subjects ▾ | Data Tools ▾ | Publications ▾ | Economic Releases ▾ | Students ▾ | Beta ▾

## Labor Force Statistics from the Current Population Survey

FONT SIZE: ⊖ ⊕

SHARE ON: f t in | CPS

**BROWSE CPS**

CPS HOME
CPS OVERVIEW
CPS NEWS RELEASES
CPS DATABASES
CPS TABLES
CPS PUBLICATIONS
CPS FAQS
CONTACT CPS

**SEARCH CPS**

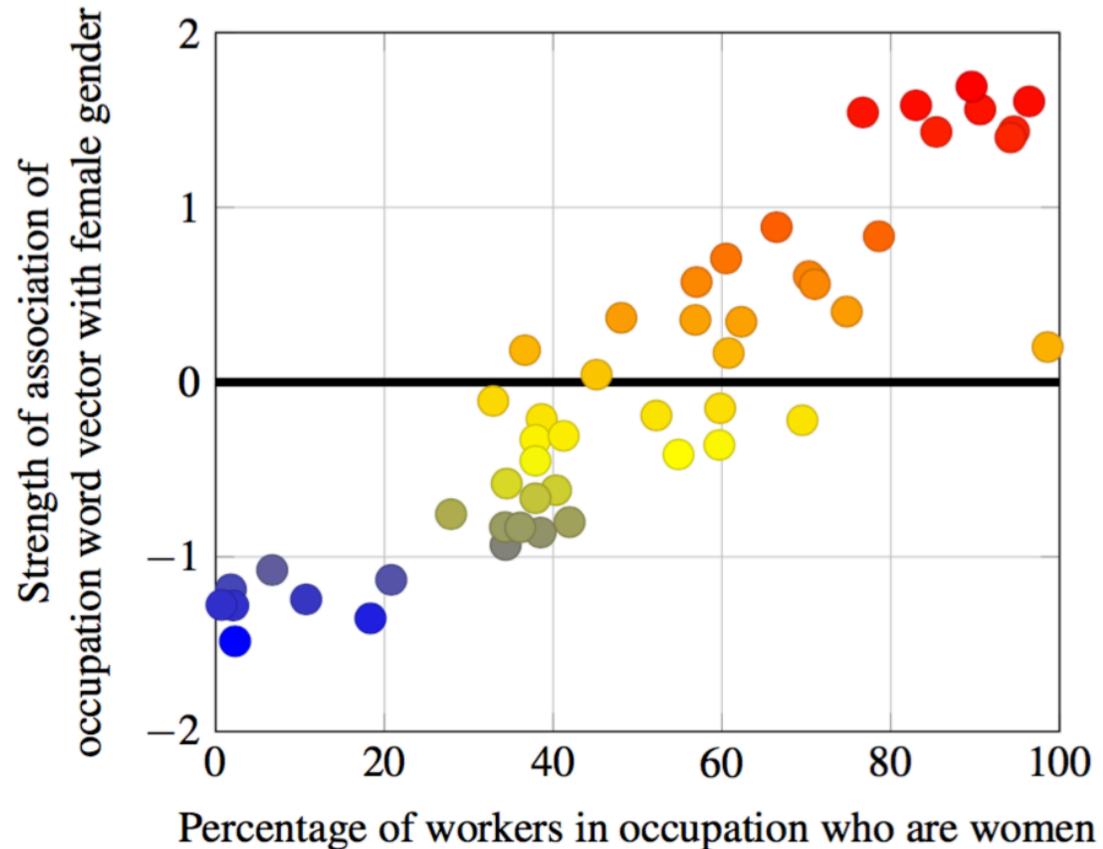[        ] Go

**CPS TOPICS**

**HOUSEHOLD DATA**
**ANNUAL AVERAGES**
**11. Employed persons by detailed occupation, sex, race, and Hispanic or Latino ethnicity**
[Numbers in thousands]

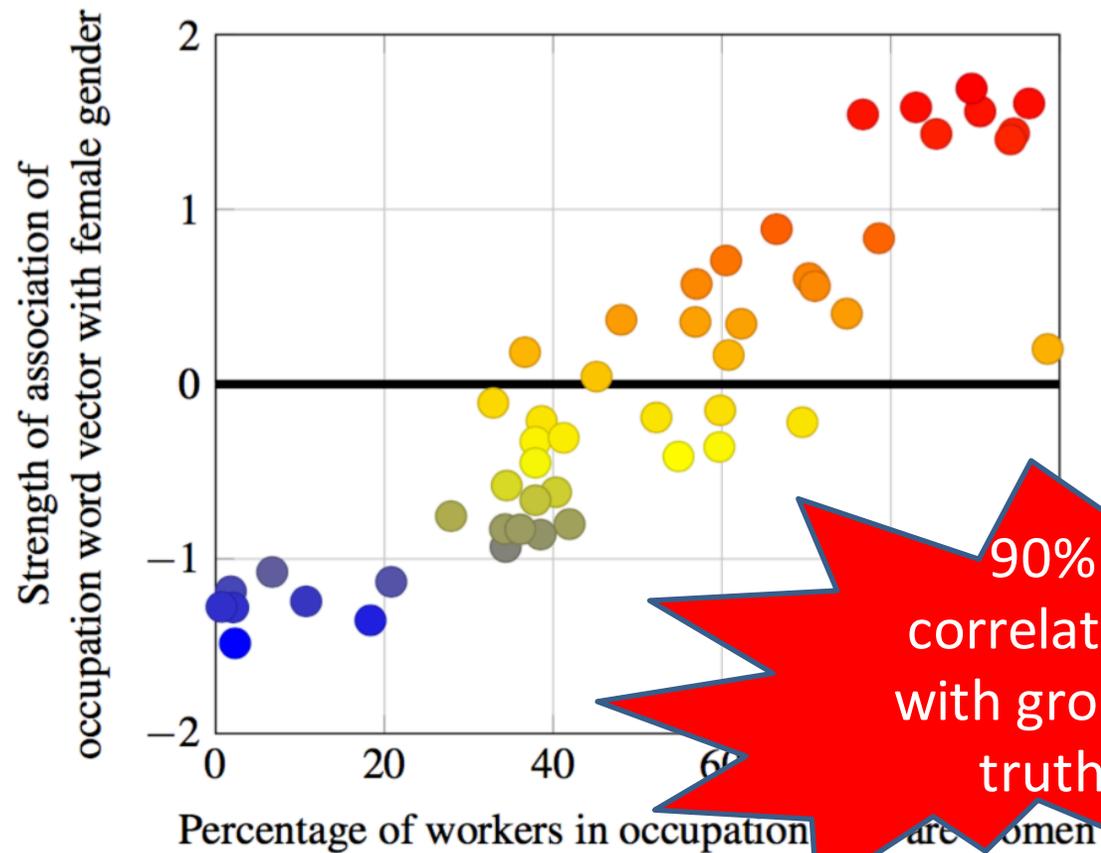| Occupation | 2015 | | | | |
|---|---|---|---|---|---|
| | | Percent of total employed | | | |
| | Total employed | Women | Black or African American | Asian | Hispanic or Latino |
| **Total, 16 years and over** | 148,834 | 46.8 | 11.7 | 5.8 | 16.4 |
| | | | | | |
| **Management, professional, and related occupations** | 57,960 | 51.5 | 9.2 | 7.7 | 9.1 |
| **Management, business, and financial operations occupations** | 24,108 | 43.6 | 8.2 | 6.3 | 9.4 |
| **Management occupations** | 16,994 | 39.2 | 7.3 | 5.6 | 9.7 |
| **Chief executives** | 1,517 | 27.9 | 3.6 | 4.7 | 5.5 |

# WEFAT: Women employed in the US



Occupation-gender association
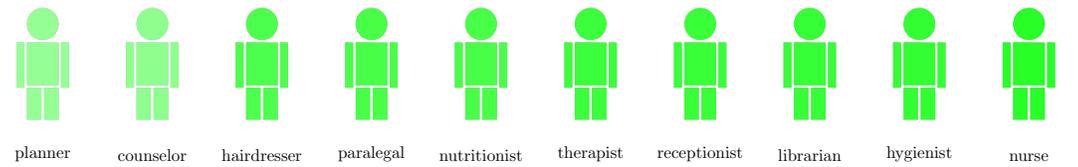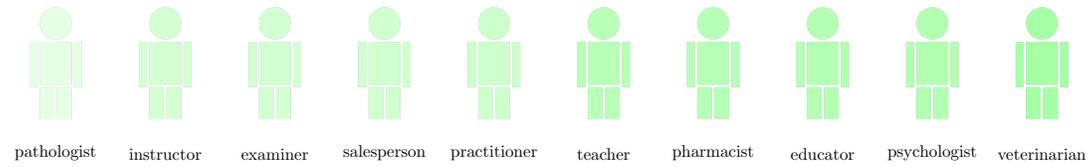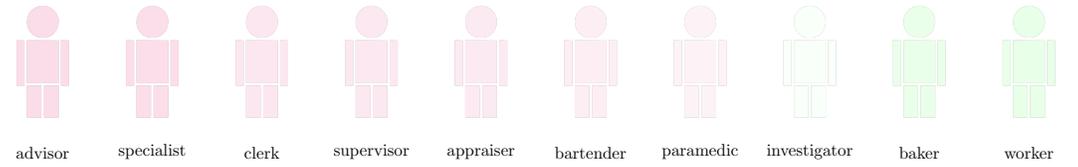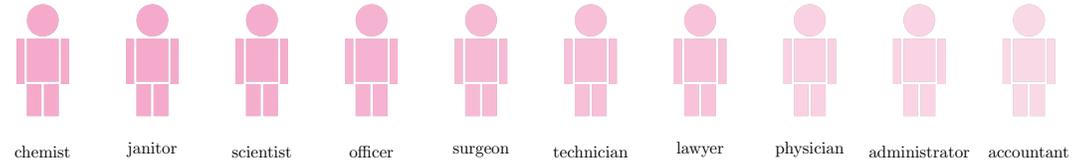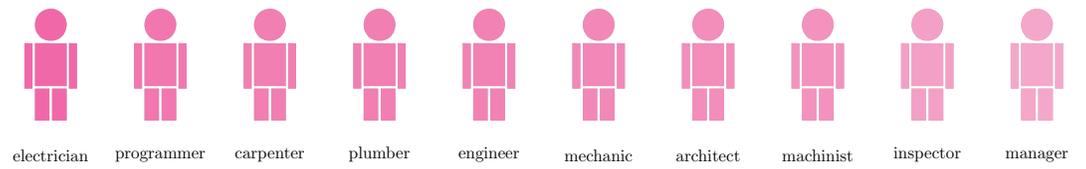Pearson's correlation coefficient $\rho = 0.90$ with $p$-value $< 10^{-18}$.

# WEFAT: Women employed in the US



Occupation-gender association
Pearson's correlation coefficient $\rho = 0.90$ with $p$-value $< 10^{-18}$.

electrician  programmer  carpenter  plumber  engineer  mechanic  architect  machinist  inspector  manager

chemist  janitor  scientist  officer  surgeon  technician  lawyer  physician  administrator  accountant

advisor  specialist  clerk  supervisor  appraiser  bartender  paramedic  investigator  baker  worker

pathologist  instructor  examiner  salesperson  practitioner  teacher  pharmacist  educator  psychologist  veterinarian

planner  counselor  hairdresser  paralegal  nutritionist  therapist  receptionist  librarian  hygienist  nurse

0            50            100

Predicted Percentage of Women with Occupation

Pearson's correlation coefficient $\rho = 0.90$ with 2015 U.S. Bureau of Labor Statistics

# Problem

# Problem

# Problem

# Problem

# Problem



O bir profesor.

# Problem

# Problem

# Problem

# True for German

# True for German

# True for Bulgarian

# True for Bulgarian



O bir doktor. → Той е лекар.

O bir hemşire. Edit → Тя е медицинска сестра.
Tya e meditsinska sestra.

# Universally Accepted Stereotypes

| Targets | Stereotype | Percentile | Effect Size |
|---|---|---|---|
| Flowers | Pleasant | | |
| Insects | Unpleasant | $10^{-8}$ | 1.35 |
| Musical Instruments | Pleasant | | |
| Weapons | Unpleasant | $10^{-7}$ | 1.53 |

**Cohen** suggested that
$|d| = 0.2$ is a 'small' **effect size**,
$|d| = 0.5$ is a 'medium' **effect size**,
$|d| >= 0.8$ is a '**large**' effect size.

# Race and Gender Stereotypes

| Targets | Stereotype | Percentile | Effect Size |
|---|---|---|---|
| White | Pleasant | $10^{-8}$ | 1.41 |
| Black | Unpleasant | | |
| Male | Career | $10^{-3}$ | 1.81 |
| Female | Family | | |
| Male | Science | $10^{-2}$ | 1.24 |
| Female | Arts | | |

**Cohen** suggested that $|d|= 0.2$ is a 'small' **effect size**, $|d|= 0.5$ is a 'medium' **effect size**, $|d|>=0.8$ is a **'large' effect size**.

# Age and Disease Stereotypes

| Targets | Stereotype | Percentile | Effect Size |
|---|---|---|---|
| Young | **Pleasant** | $10^{-2}$ | 1.21 |
| Old | **Unpleasant** | | |
| **Physical Disease** | **Controllable** | $10^{-2}$ | 1.67 |
| **Mental Disease** | **Uncontrollable** | | |

**Cohen** suggested that $|d| = 0.2$ is a 'small' **effect size**, $|d| = 0.5$ is a 'medium' **effect size**, $|d| >= 0.8$ is a **'large' effect size**.

# Sexual Stigma and Transphobia

| Targets | Stereotype | Percentile | Effect Size |
|---------|-----------|:----------:|:-----------:|
| Heterosexual | Pleasant | $10^{-2}$ | 1.27 |
| Homosexual | Unpleasant | | |
| Straight | Pleasant | $10^{-2}$ | 1.34 |
| Transgender | Unpleasant | | |

**Cohen** suggested that
$|d|$= 0.2 is a 'small' **effect size**,
$|d|$= 0.5 is a 'medium' **effect size**,
$|d|>=0.8$ is a **'large' effect size**.

# German: Gender Stereotypes and Nationalism

| Targets | Stereotype | Percentile | Effect Size |
|---------|-----------|------------|-------------|
| Male | Career | $10^{-2}$ | 1.54 |
| Female | Family | | |
| Male | Science | $10^{-2}$ | 1.56 |
| Female | Arts | | |
| German | Pleasant | $10^{-2}$ | 1.34 |
| Turkish | Unpleasant | | |

**Cohen** suggested that
$|d| = 0.2$ is a 'small' **effect size**,
$|d| = 0.5$ is a 'medium' **effect size**,
$|d| >= 0.8$ is a **'large' effect size**.
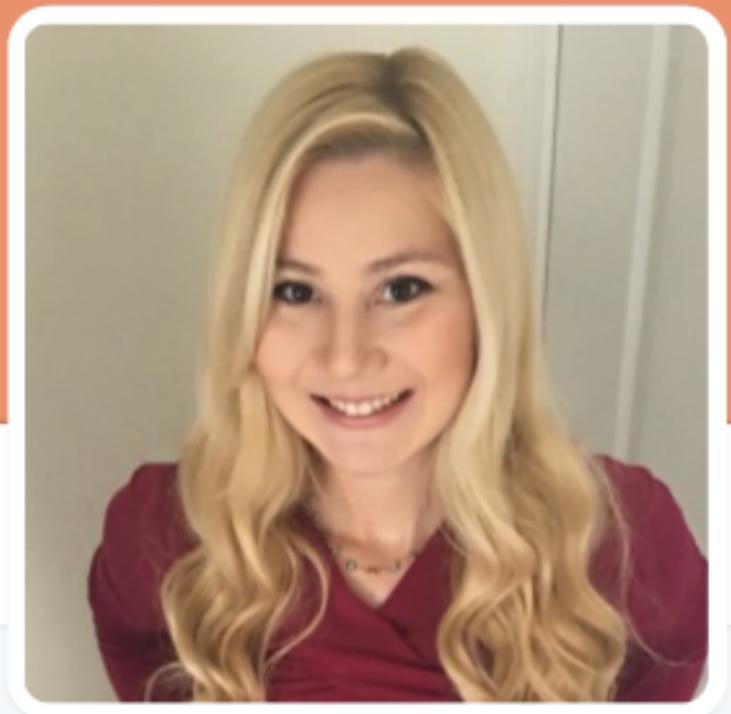
# Discussion points:

- Machine learning expertise for algorithmic transparency
- How to mitigate bias while preserving utility
- How long does bias persist in models?
- Are biased models causing a snowball effect?
- Policy to stop discrimination
  - predictive policing
  - ML services effect billions every day
    - Google, Amazon, Microsoft

## Research Code

github.com/calaylin

## Webpage

princeton.edu/~aylinc

## Check our blog

freedom-to-tinker.com

**Aylin Caliskan**
@aylin_cim

FREEDOM TO TINKER
research and expert commentary on digital technologies in public life