# ALGORITHMIC MIRRORS OF SOCIETY

## Aylin Caliskan
## @aylin_cim

Assistant Professor
Department of Computer Science
The George Washington University

# Garbage in garbage out



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

# Garbage in garbage out

# Problem

# Problem

# Problem

# Problem

# Account of bias in AI

human bias → semantics → distributional meaning → ML models

Semantics derived automatically from language corpora contain human-like biases

Aylin Caliskan[1,*], Joanna J. Bryson[1,2,*], Arvind Narayanan[1,*]

+ See all authors and affiliations

# Account of bias in AI

# Natural language processing as a service:

# Language models used by billions every day:

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

# Generating language models

# Word embeddings

- Syntax
- Analogies
- Semantic similarity

$\text{word}_1, \text{feature}_1, \text{feature}_2, \qquad \text{feature}_{300}$
$\text{word}_2, \text{feature}_1, \text{feature}_2, \qquad \text{feature}_{300}$
$\text{word}_3, \text{feature}_1, \text{feature}_2, \qquad \text{feature}_{300}$
$\qquad \dots \quad \dots$
$\text{word}_{2,000,000}, \text{feature}_1, \text{feature}_2, \text{feature}_{300}$

# Vector arithmetic

# Universally Accepted Stereotypes

| Targets | Stereotype | Percentile | Effect Size |
|---|---|---|---|
| Flowers | Pleasant | $10^{-7}$ | 1.50 |
| Insects | Unpleasant | | |
| Musical Instruments | Pleasant | $10^{-7}$ | 1.53 |
| Weapons | Unpleasant | | |

**Statistically significant
&
Large effect size (Cohen's d)**

# 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 | | |

**Statistically significant**
**&**
**Large effect size (Cohen's d)**

# 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 | | | | |
| --- | --- | --- | --- | --- | --- |
| | Total employed | Percent of 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
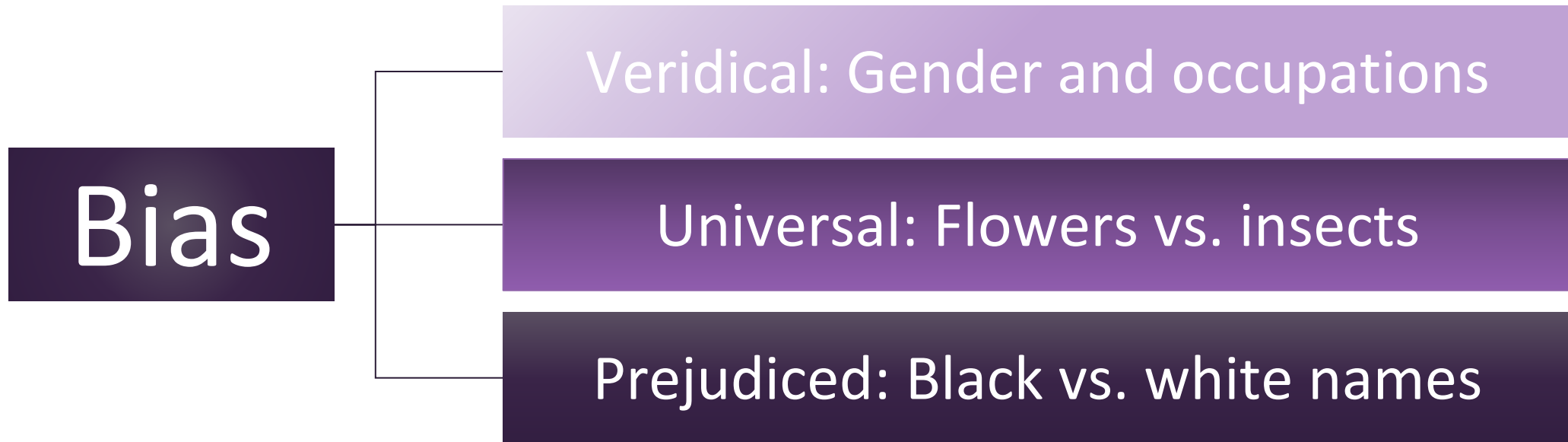Pearson's correlation coefficient $\rho = 0.90$ with $p$-value $< 10^{-18}$.

# Societal bias in AI

Bias
- Veridical: Gender and occupations
- Universal: Flowers vs. insects
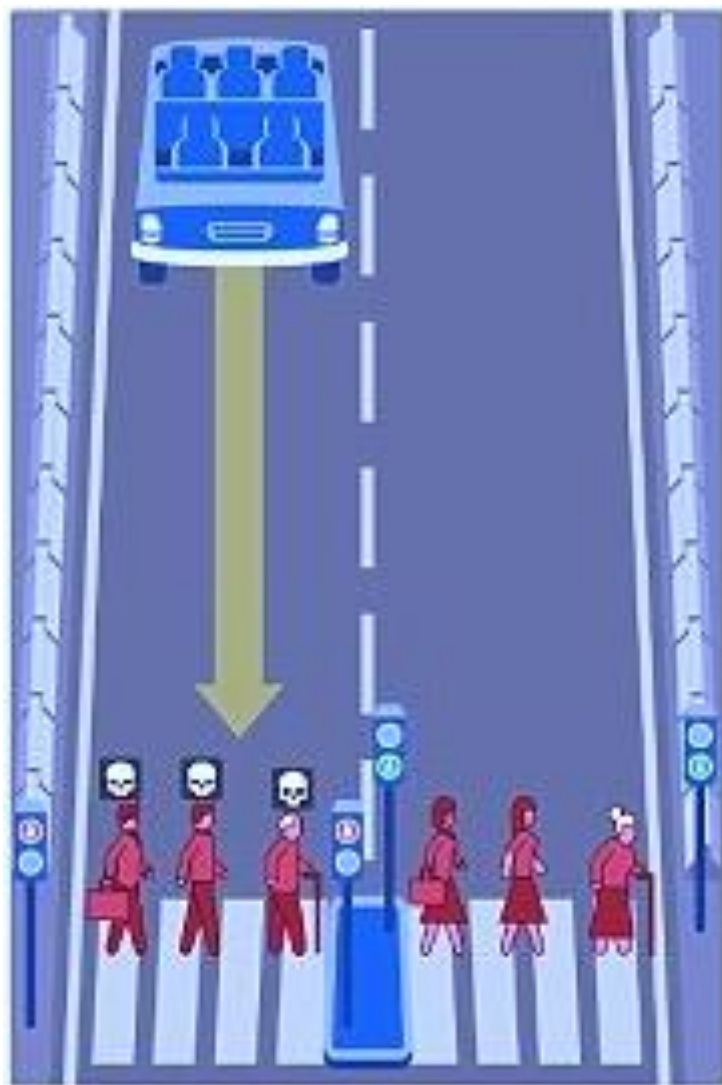- Prejudiced: Black vs. white names

- Veridical associations may result in bias/prejudice
- Prejudice := unacceptable bias (cultural, evolving)
- Same bias may be desirable or unacceptable
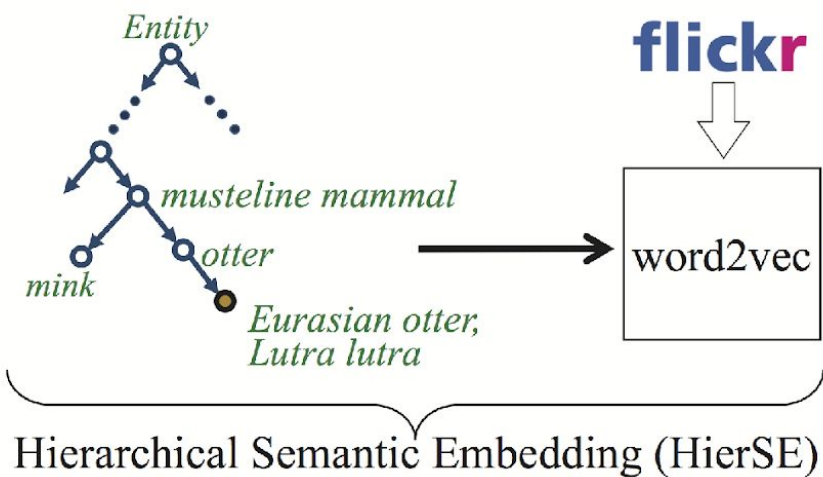
# Surveillance

# Driverless cars

Natural Language Processing

**Computer Vision**

Bias in Computer Vision

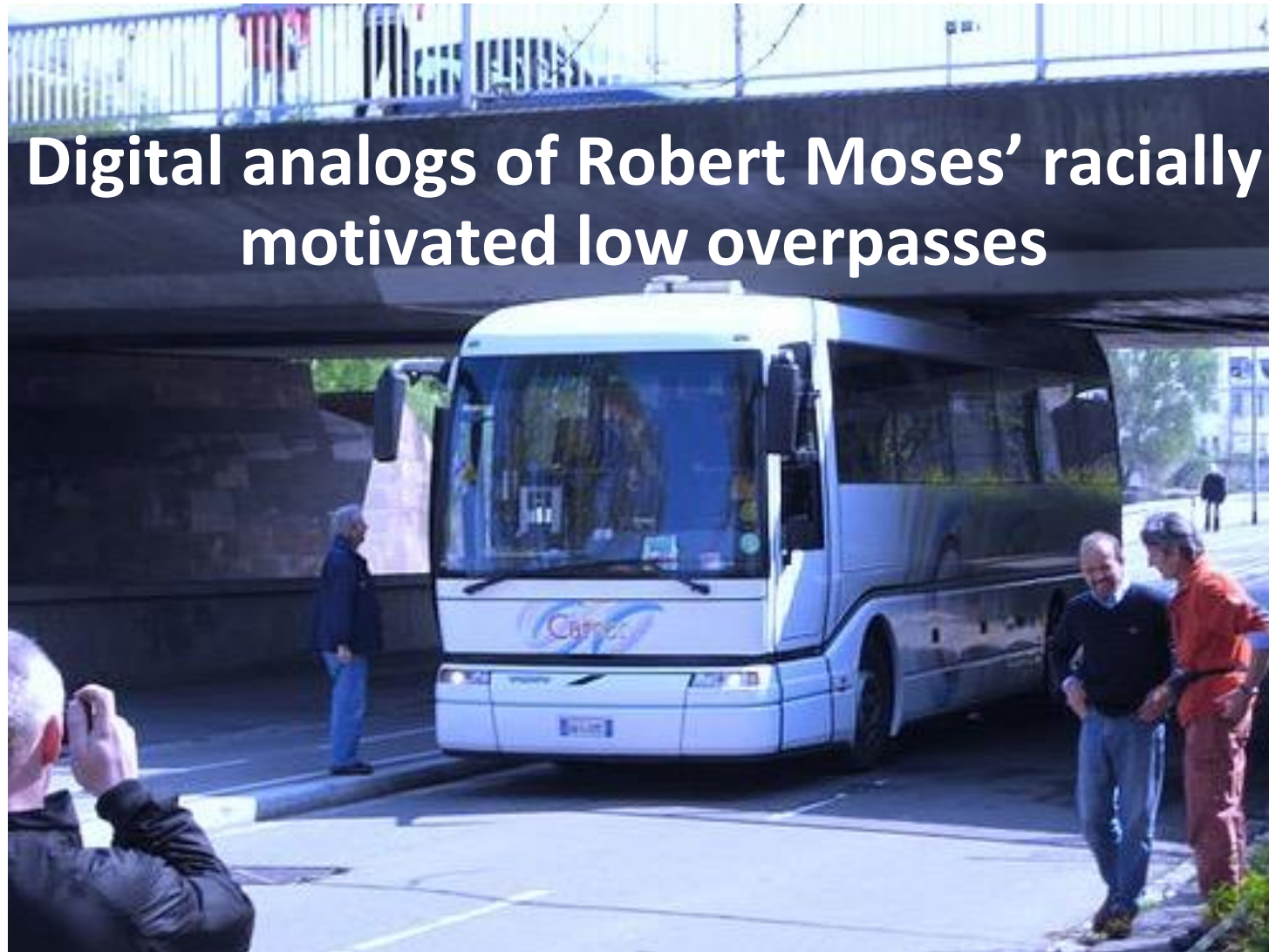| Good | Enjoy, Smiling, Cheer, Fantastic, Celebrate, Triumph, Friendship, Joyous |
| --- | --- |
| Bad | Negative, Despise, Scorn, Disgust, Yucky, Hatred, Humiliate, Sickening |
| Black people | |
| White people | |

p-value: 0.05  ---  effect size: 0.98 = **SIGNIFICANT HIGH BIAS**

# Effects of machine bias on society and AI



Digital analogs of Robert Moses' racially motivated low overpasses

# Semantics derived automatically from language corpora contain human-like biases

**Aylin Caliskan**[1,*], **Joanna J. Bryson**[1,2,*], **Arvind Narayanan**[1,*]

**+** See all authors and affiliations