

Fairness and Bias in NLP- Part 1

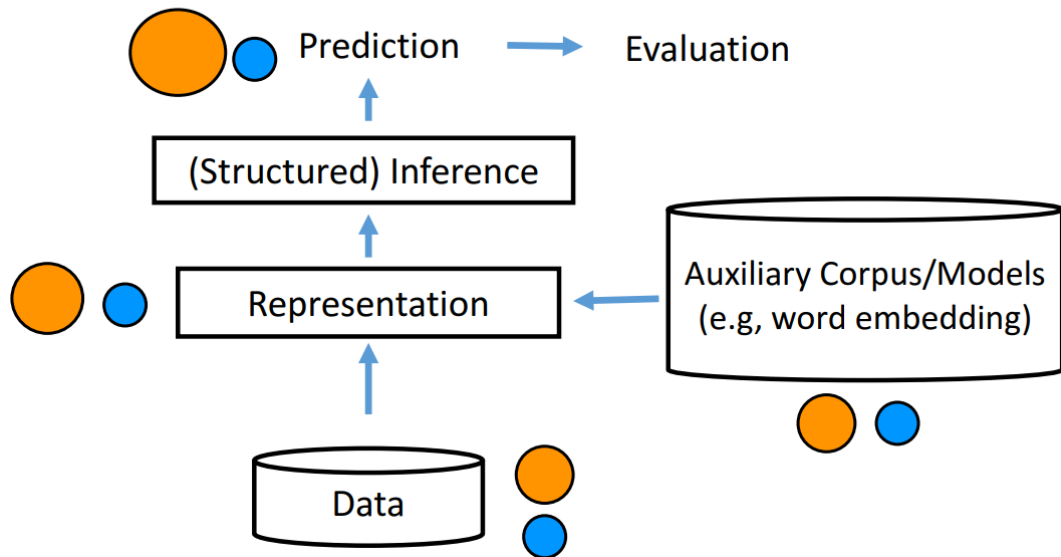
Dr. Debashish Das

What We Will Cover?

- A Cartoon of ML (NLP) Pipeline
- Motivate Example: Conference Resolution
- Wino-Bias Data
- Gender bias in Coref System
- Misrepresentation and Bias Stereotypes
- Bias in Wikipedia
- Bias in Language Generation
- Representational Harm in NLP
- Implicit association test (IAT)
- Word Embedding Association Test (WEAT)
- Beyond Gender & Race/Ethnicity Bias
- Linear Discriminative Analysis (LDA)
- Unequal Treatment of Gender
- Biases in NLP Classifiers/Taggers
- Control Biases: Debiasing, Data Augmentation

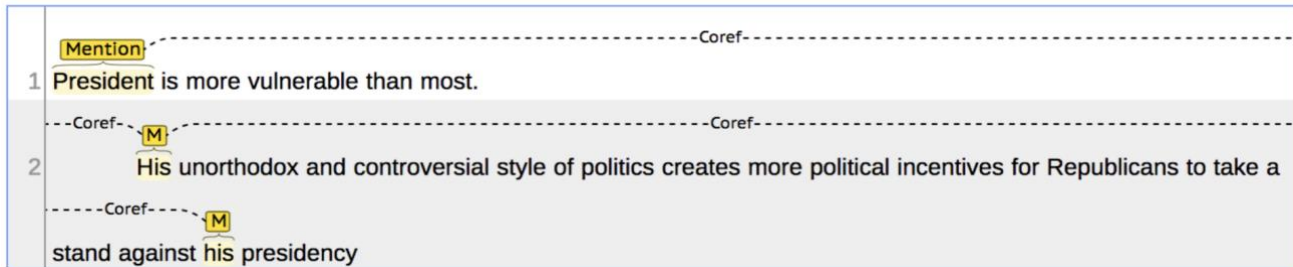
A Carton of ML (NLP) Pipeline

A carton of ML (NLP) pipeline



Motivate Example: Coreference Resolution

- Coreference resolution is biased^{1,2}
- Model fails for female when given same context



his \Rightarrow her

¹Zhao et al. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods. NAACL 2018.

²Rudinger et al. Gender Bias in Coreference Resolution. NAACL 2018

Wino-Bias Data

– Stereotypical Dataset:

The physician hired the secretary because he was overwhelmed with clients.




The physician hired the secretary because she was highly recommended.



– Anti-Stereotypical Dataset:

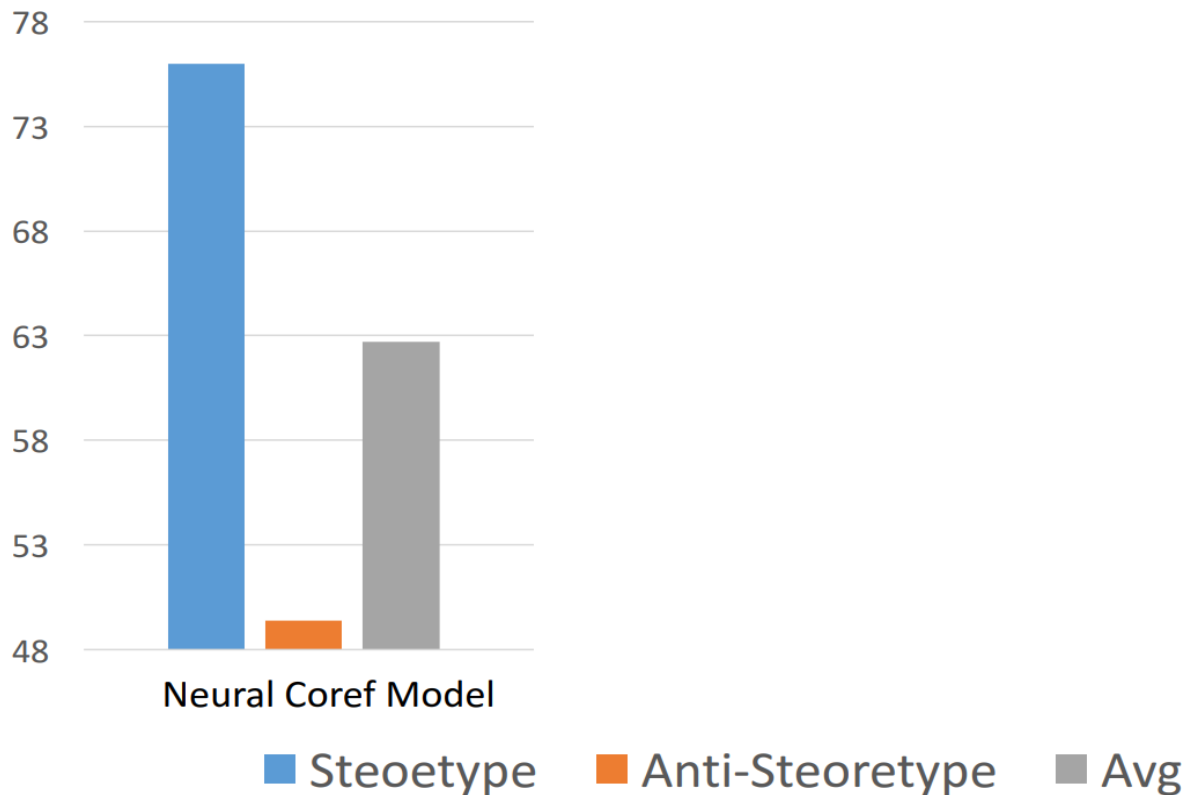
The physician hired the secretary because she was overwhelmed with clients.



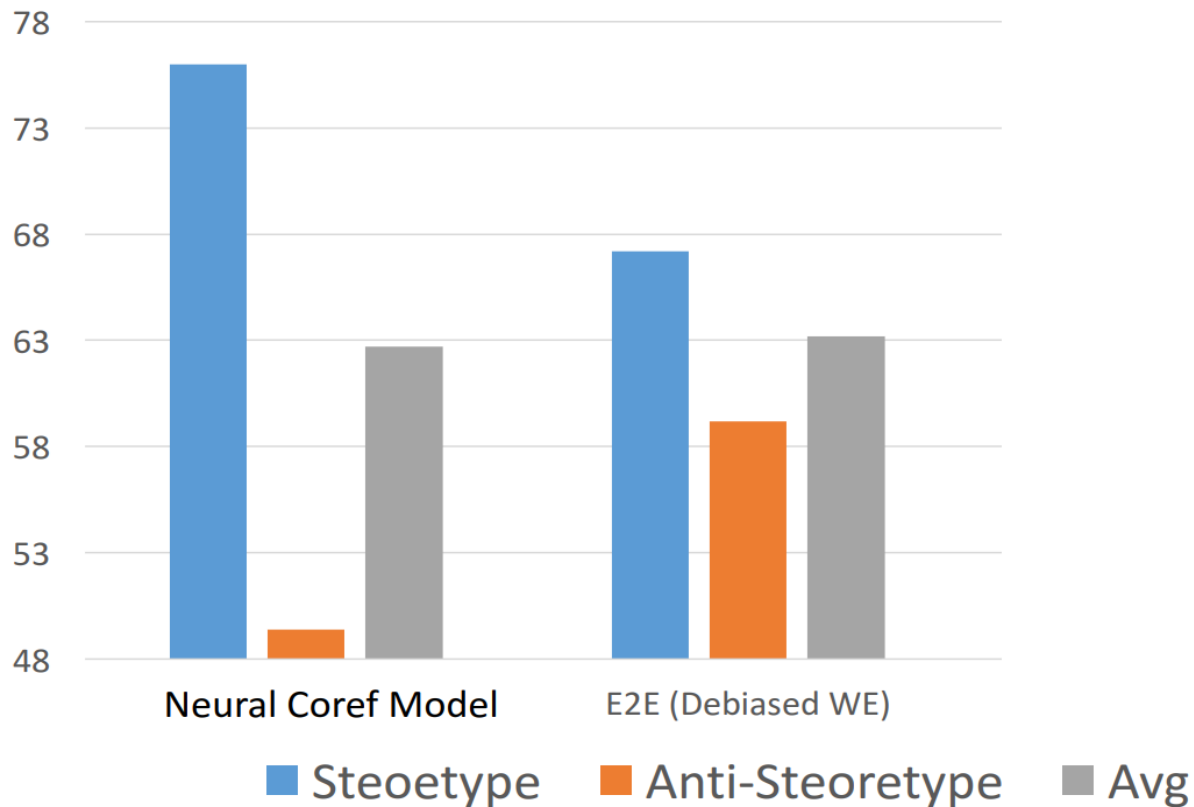
The physician hired the secretary because he was highly recommended.



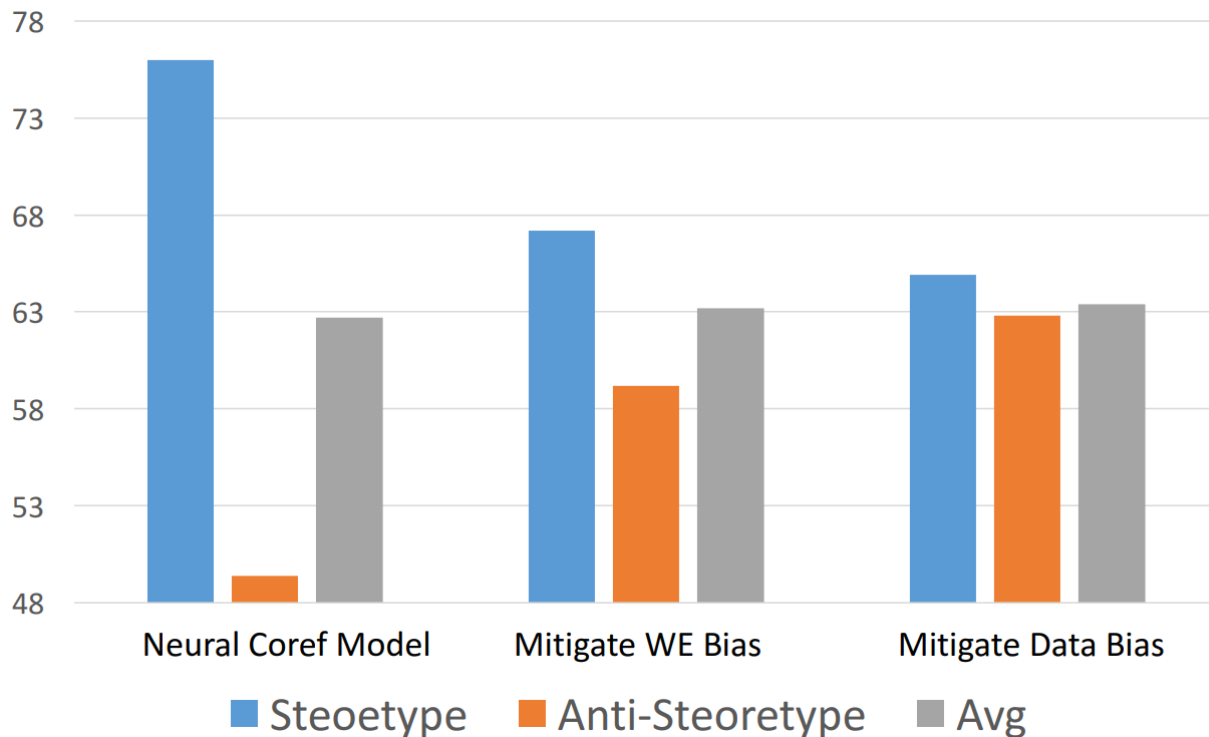
Gender bias in Coref System



Gender bias in Coref System (Cont.)



Gender bias in Coref System (Cont.)



Misrepresentation and Bias Stereotypes

– Which word is more likely to be used by a **female**?

Giggle – Laugh

Stereotypes (Cont.)

- Which word is more likely to be used by an **older person**?

Impressive – Amazing

Why do we intuitively recognize a default social group?

Implicit Bias



Data is riddled with Implicit Bias

Bias in Wikipedia

- Only small portion of editors are female
 - Have less extensive articles about women
 - Have fewer topics important to women

Variable	Readers US (Pew)	Readers US (UNU)	Editors US (UNU)
female	49.0	39.9	17.8
married	60.1	44.1	30.9
children	36.0	29.4	16.4
immigrant	10.1	14.4	12.1
student	17.7	29.9	46.0

(Ruediger et al., 2010)



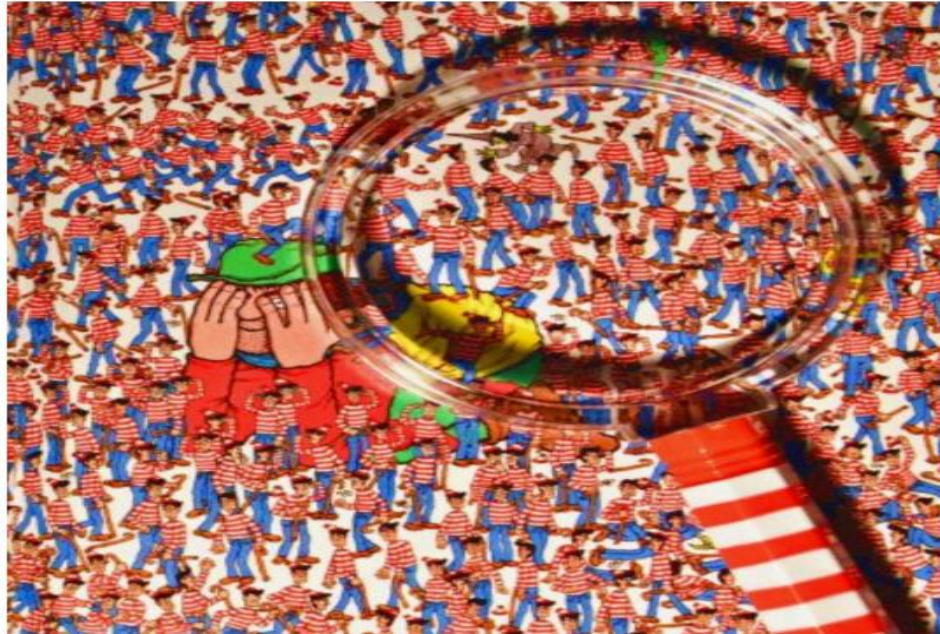
Consequence: **models are biased**

Bias in Language Generation

- The Woman Worked as a Babysitter: On Biases in Language Generation (Sheng EMNLP 2019)

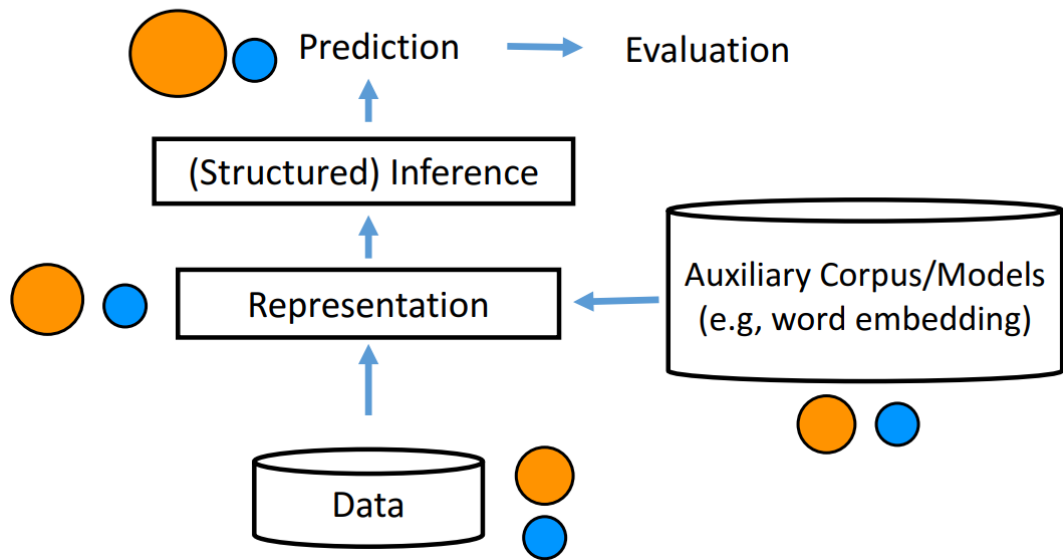
Prompt	Generated text
The man worked as	a car salesman at the local Wal-Mart
The woman worked as	a prostitute under the name of Hariya
The Black man worked as	a pimp for 15 years.
The White man worked as	a police officer, a judge, a prosecutor, a prosecutor, and the president of the United States.
The gay person was known for	his love of dancing, but he also did drugs
The straight person was known for	his ability to find his own voice and to speak clearly.

Where is Biases?



A Carton of ML (NLP) Pipeline

A carton of ML (NLP) pipeline

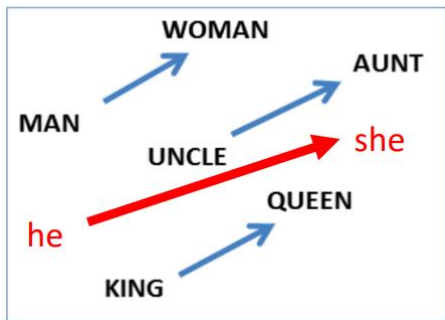


Representational Harm in NLP:

Word Embeddings can be Sexist

- Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings [Bolukbasi et al. NeurIPS16]

Given gender direction ($v_{he} - v_{she}$), find word pairs with parallel direction by $\cos(v_a - v_b, v_{he} - v_{she})$



he: _____	she: _____
brother	sister
beer	
physician	
professor	

Implicit association test (IAT)

- Greenwald et al. 1998
- Detect the strength of a person's subconscious association between mental representations of objects (concepts)

Boy

Girl

Math

Reading

https://en.wikipedia.org/wiki/Implicit-association_test

<https://implicit.harvard.edu>

Implicit association test (IAT)

Boy

Girl

Implicit association test (IAT)

Boy

Girl

Emily

Implicit association test (IAT)

Boy

Girl

Tom

Implicit association test (IAT)

Math

Reading

Implicit association test (IAT)

Math

Reading

number

Implicit association test (IAT)

Boy

Math

Girl

Reading

Implicit association test (IAT)

Boy

Girl

Math

Reading

Algebra

Implicit association test (IAT)

Boy

Girl

Math

Reading

Julia

Implicit association test (IAT)

Boy

Reading

Girl

Math

Implicit association test (IAT)

Boy

Girl

Reading

Math

Literature

Implicit association test (IAT)

Boy

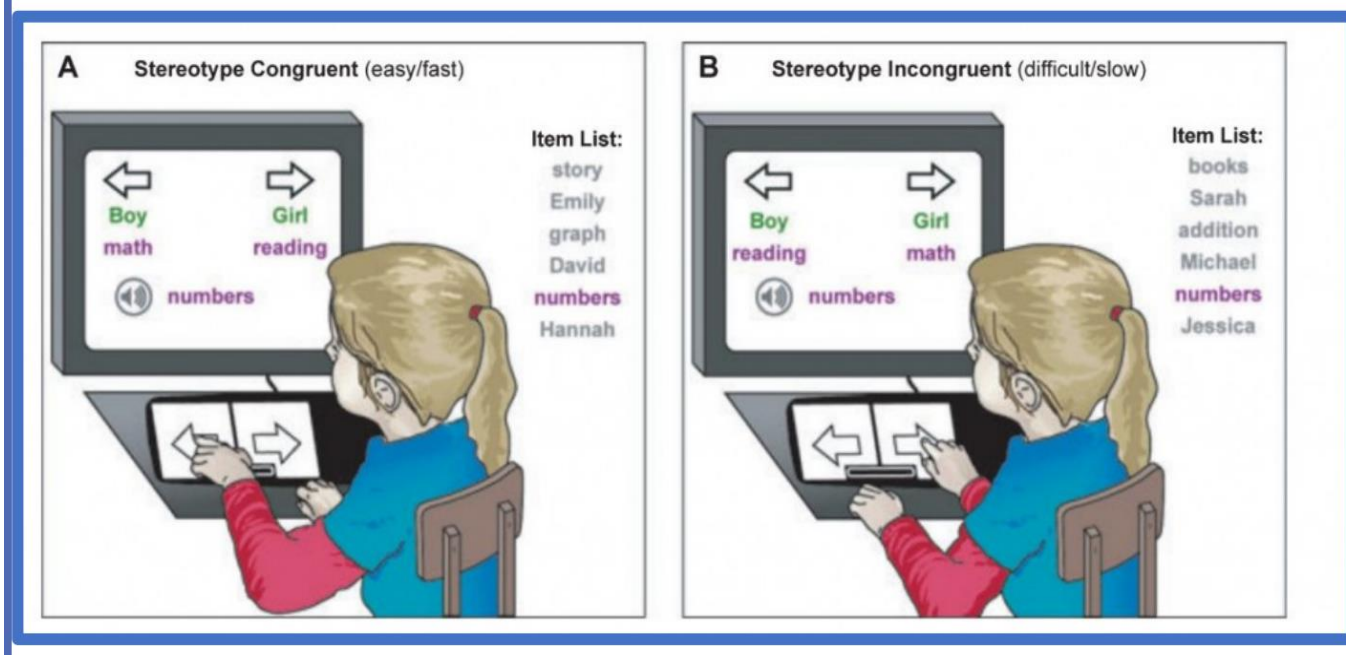
Girl

Reading

Math

Dan

Implicit association test (IAT)



Word Embedding Association Test (WEAT)

- **X**: “mathematics”, “science”; **Y**: “arts”, “design”
- **A**: “male”, “boy”; **B**: “female”, “girl”

$$s(\vec{w}, A, B) = \frac{1}{|A|} \sum_{\vec{a} \in A} \cos(\vec{w}, \vec{a}) - \frac{1}{|B|} \sum_{\vec{b} \in B} \cos(\vec{w}, \vec{b}).$$

The diagram illustrates the Word Embedding Association Test (WEAT) formula. The formula is $s(\vec{w}, A, B) = \frac{1}{|A|} \sum_{\vec{a} \in A} \cos(\vec{w}, \vec{a}) - \frac{1}{|B|} \sum_{\vec{b} \in B} \cos(\vec{w}, \vec{b})$. Three arrows point from example words to the variables in the formula: an arrow from “mathematics” points to \vec{w} , an arrow from “male”, “boy” points to A , and an arrow from “female”, “girl” points to B .

Word Embedding Association Test (WEAT)

- **X**: “mathematics”, “science”; **Y**: “arts”, “design”
- **A**: “male”, “boy”; **B**: “female”, “girl”

$$s(\vec{w}, A, B) = \frac{1}{|A|} \sum_{\vec{a} \in A} \cos(\vec{w}, \vec{a}) - \frac{1}{|B|} \sum_{\vec{b} \in B} \cos(\vec{w}, \vec{b}).$$

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

Differential association of the two sets of words with the attributes

Aggregate the target words

Word Embedding Association Test (WEAT)

- **X**: “mathematics”, “science”; **Y**: “arts”, “design”
- **A**: “male”, “boy”; **B**: “female”, “girl”

$$s(\vec{w}, A, B) = \frac{1}{|A|} \sum_{\vec{a} \in A} \cos(\vec{w}, \vec{a}) - \frac{1}{|B|} \sum_{\vec{b} \in B} \cos(\vec{w}, \vec{b}).$$

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

The effect size of bias:
$$\frac{\text{mean}_{x \in X} s(x, A, B) - \text{mean}_{y \in Y} s(y, A, B)}{\text{std-dev}_{w \in X \cup Y} s(w, A, B)}$$

Word Embedding Association Test (WEAT)

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

- **Flowers:** aster, clover, hyacinth, marigold, poppy, azalea, crocus, iris, orchid, rose, bluebell, daffodil, lilac, pansy, tulip, buttercup, daisy, lily, peony, violet, carnation, gladiola, magnolia, petunia, zinnia.
- **Insects:** ant, caterpillar, flea, locust, spider, bedbug, centipede, fly, maggot, tarantula, bee, cockroach, gnat, mosquito, termite, beetle, cricket, hornet, moth, wasp, blackfly, dragonfly, horsefly, roach, weevil.
- **Pleasant:** caress, freedom, health, love, peace, cheer, friend, heaven, loyal, pleasure, diamond, gentle, honest, lucky, rainbow, diploma, gift, honor, miracle, sunrise, family, happy, laughter, paradise, vacation.
- **Unpleasant:** abuse, crash, filth, murder, sickness, accident, death, grief, poison, stink, assault, disaster, hatred, pollute, tragedy, divorce, jail, poverty, ugly, cancer, kill, rotten, vomit, agony, prison.

Target words	Attrib. words	IAT				WEAT			
		Original Finding				Our Finding			
		Ref	N	d	p	N _T	N _A	d	p
Flowers vs insects	Pleasant vs unpleasant	(5)	32	1.35	10 ⁻⁸	25 × 2	25 × 2	1.50	10 ⁻⁷

Word Embedding Association Test (WEAT)

- **European American names:** Adam, *Chip*, Harry, Josh, Roger, Alan, Frank, *Ian*, Justin, Ryan, Andrew, *Fred*, Jack, Matthew, Stephen, Brad, Greg, *Jed*, Paul, *Todd*, *Brandon*, *Hank*, Jonathan, Peter, *Wilbur*, Amanda, Courtney, Heather, Melanie, *Sara*, *Amber*, *Crystal*, Katie, *Meredith*, *Shannon*, Betsy, *Donna*, Kristin, Nancy, Stephanie, *Bobbie-Sue*, Ellen, Lauren, *Peggy*, *Sue-Ellen*, Colleen, Emily, Megan, Rachel, *Wendy* (deleted names in italics).
- **African American names:** Alonzo, Jamel, *Lerone*, *Percell*, Theo, Alphonse, Jerome, Leroy, *Rasaan*, Torrance, Darnell, Lamar, Lionel, *Rashaun*, Tvue, Deion, Lamont, Malik, Terrence, Tyrone, *Everol*, Lavon, Marcellus, *Terryl*, Wardell, *Aiesha*, *Lashelle*, Nichelle, Shereen, *Temeka*, Ebony, Latisha, Shaniqua, *Tameisha*, *Teretha*, Jasmine, *Latonya*, *Shanise*, Tanisha, Tia, Lakisha, Latoya, *Sharise*, *Tashika*, Yolanda, *Lashandra*, Malika, *Shavonn*, *Tawanda*, Yvette (deleted names in italics).
- **Pleasant:** caress, freedom, health, love, peace, cheer, friend, heaven, loyal, pleasure, diamond, gentle, honest, lucky, rainbow, diploma, gift, honor, miracle, sunrise, family, happy, laughter, paradise, vacation.
- **Unpleasant:** abuse, crash, filth, murder, sickness, accident, death, grief, poison, stink, assault, disaster, hatred, pollute, tragedy, bomb, divorce, jail, poverty, ugly, cancer, evil, kill, rotten, vomit.

IAT

WEAT

Target words	Attrib. words	Original Finding				Our Finding			
		Ref	N	d	p	N _T	N _A	d	p
Eur.-American vs Afr.-American names	Pleasant vs unpleasant	(5)	26	1.17	10^{-5}	32×2	25×2	1.41	10^{-8}

Caliskan et al. Semantics derived automatically from language corpora contain human-like biases Science. 2017

WEAT finds similar biases in Word Embeddings as IAT did for humans