



UNIVERSITY OF
GOTHENBURG

SPRÅKBANKEN **TEXT**

Intrinsic Bias Metrics Do Not Correlate with Application Bias

Ricardo Muñoz Sánchez

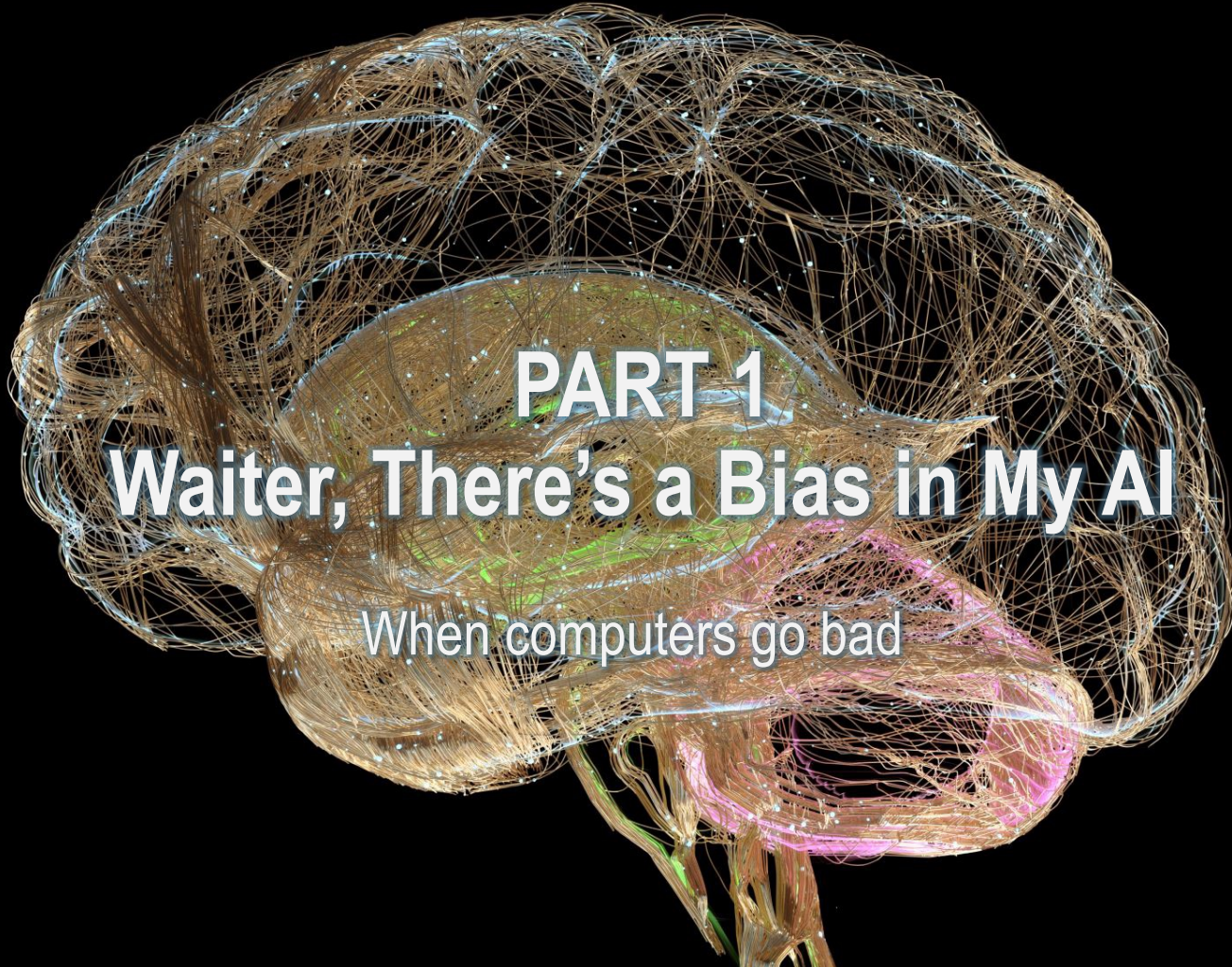
Based on work done with:

Seraphina Goldfarb-Tarrant, Rebecca Marchant, Mugdha Pandya, and Adam Lopez

Overview

- Biases in AI
- How do we measure them?
- Can they be removed?
- Even if we can, does it do *anything* at all?





PART 1

Waiter, There's a Bias in My AI

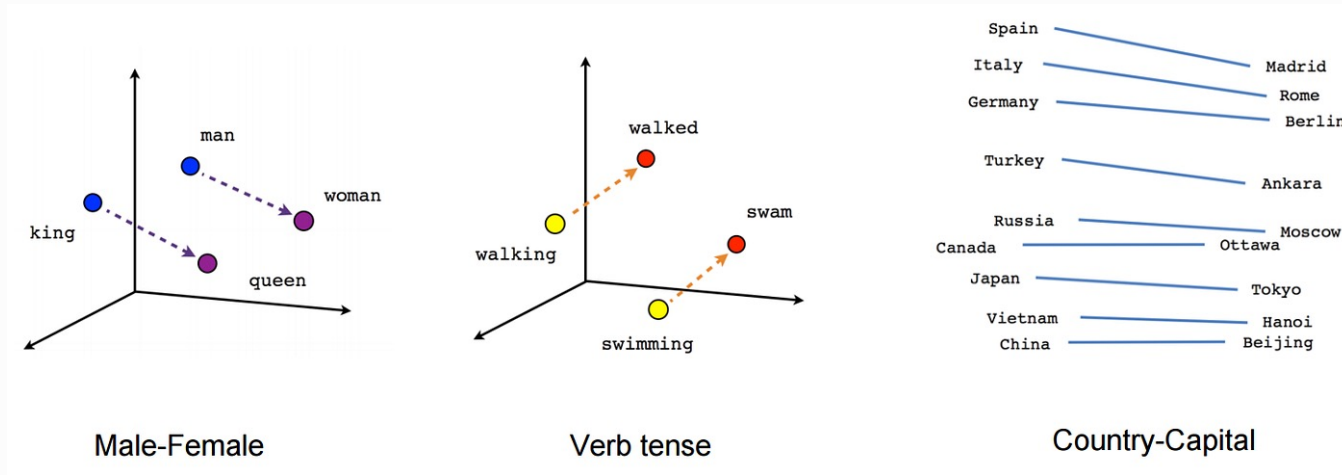
When computers go bad

Word Embeddings

- Ways to represent words in a computer-interpretable manner
- They encode both semantics and syntax of words
- Have nice geometric properties



Word Embeddings – Geometry



Word Embeddings – Analogies

- *Man is to king as woman is to... queen*
- *Walk is to swim as walking is to... swimming*
- *Spain is to Madrid as Italy is to... Rome*

Word Embeddings – Analogies

- *Man is to king as woman is to... queen*
- *Walk is to swim as walking is to... swimming*
- *Spain is to Madrid as Italy is to... Rome*
- *Man is to programmer as woman is to...*

Word Embeddings – Analogies

- *Man is to king as woman is to... queen*
- *Walk is to swim as walking is to... swimming*
- *Spain is to Madrid as Italy is to... Rome*
- *Man is to programmer as woman is to... homemaker*

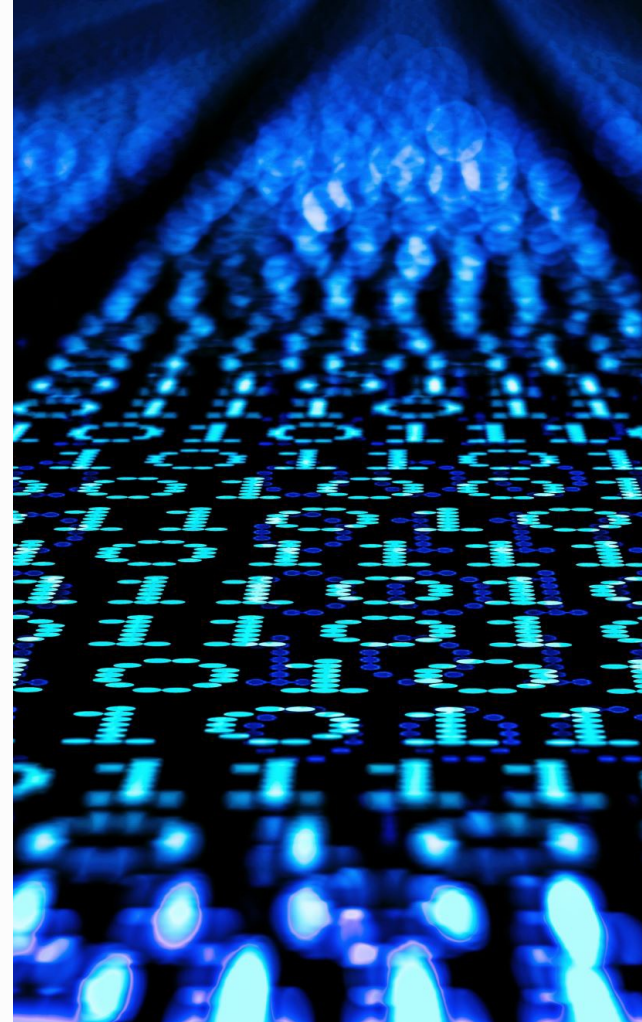
Wait, what?

Encoding Biases

- As with all AI models, embeddings find and exploit patterns in the data
- However, stereotypes are patterns in the data
- Our systems can and will pick these patterns and perpetuate unwanted biases, such as sexism and racism

Encoding Biases

- As with all AI models, embeddings find and exploit patterns in the data
- However, humans are biased and this is reflected in the data we produce
- Our models can and will pick these patterns and perpetuate unwanted biases



Correference Resolution

- The **doctor** hired a nurse because **he** was busy (**Correct**)



Correference Resolution

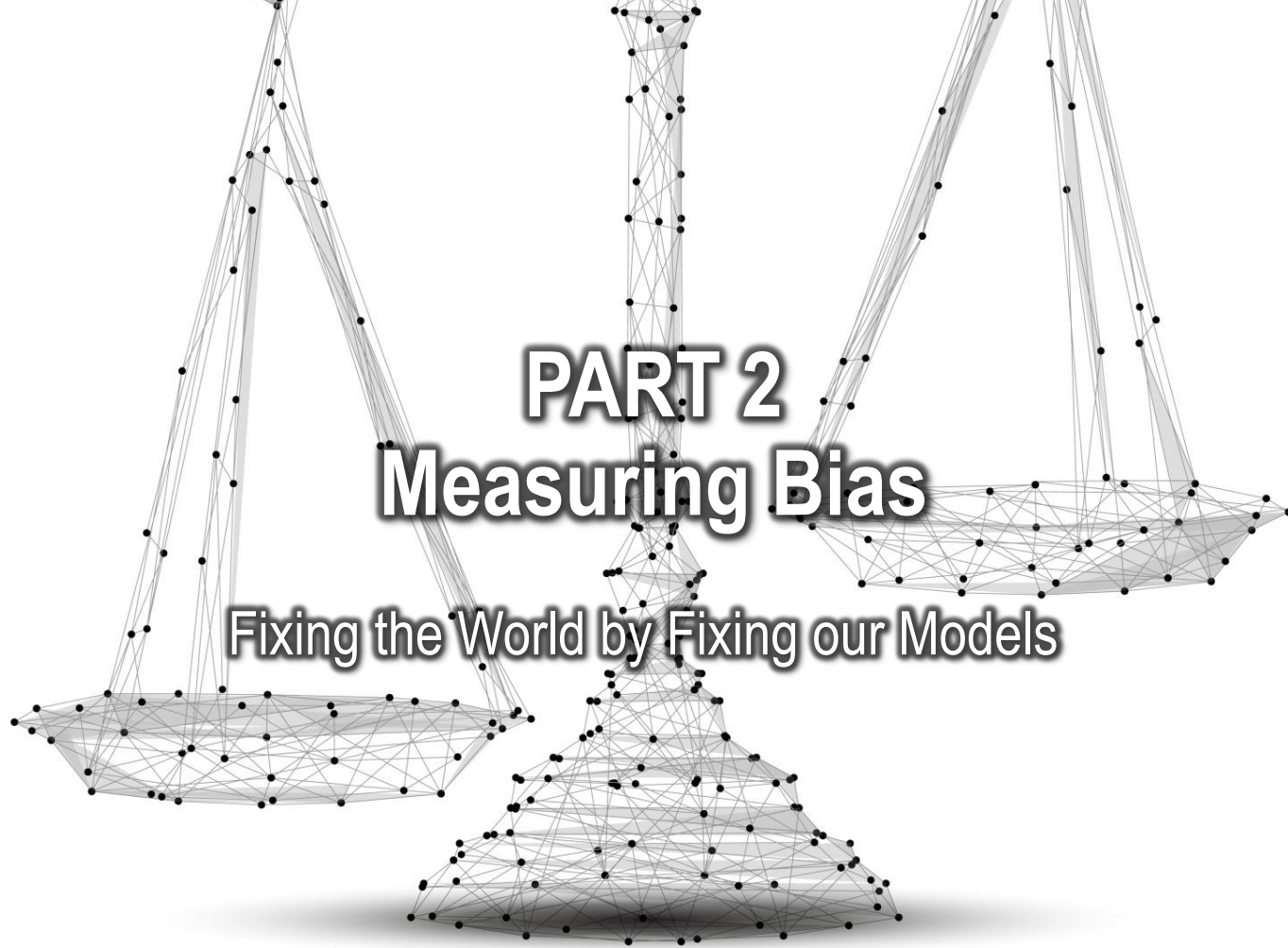
- The **doctor** hired a nurse because **he** was busy (**Correct**)
- The doctor hired a **nurse** because **she** was busy (**Wrong**)



Correference Resolution

- The **doctor** hired a nurse because **he** was busy (**Correct**)
- The doctor hired a **nurse** because **she** was busy (**Wrong**)
- The **doctor** hired a nurse because **she** was busy (**Correct**)





PART 2

Measuring Bias

Fixing the World by Fixing our Models

Bias Metrics

Intrinsic

- Measure unwanted associations in language models and embeddings
- Examples are WEAT, CEAT, and DisCo
- Also called bias metrics

Extrinsic

- Measure the disparity of performance in downstream applications
- Examples are demographic parity, equality of opportunity, and predictive rate parity
- Also called fairness metrics

Intrinsic Metrics - WEAT



Based on the Implicit Association Test (IAT)



We are given a set of target words and two lists of characteristics (stereotypical and anti-stereotypical)



Are the targets' representations closer to their stereotypes'?

Intrinsic Metrics - WEAT

Test	Target Set #1	Target Set #2	Attribute Set #1	Attribute Set #2
T1	Flowers (e.g., <i>aster</i> , <i>tulip</i>)	Insects (e.g., <i>ant</i> , <i>flea</i>)	Pleasant (e.g., <i>health</i> , <i>love</i>)	Unpleasant (e.g., <i>abuse</i>)
T2	Instruments (e.g., <i>cello</i> , <i>guitar</i>)	Weapons (e.g., <i>gun</i> , <i>sword</i>)	Pleasant	Unpleasant
T3	Euro-American names (e.g., <i>Adam</i>)	Afro-American names (e.g., <i>Jamel</i>)	Pleasant (e.g., <i>caress</i>)	Unpleasant (e.g., <i>abuse</i>)
T4	Euro-American names (e.g., <i>Brad</i>)	Afro-American names (e.g., <i>Hakim</i>)	Pleasant	Unpleasant
T5	Euro-American names	Afro-American names	Pleasant (e.g., <i>joy</i>)	Unpleasant (e.g., <i>agony</i>)
T6	Male names (e.g., <i>John</i>)	Female names (e.g., <i>Lisa</i>)	Career (e.g. <i>management</i>)	Family (e.g., <i>children</i>)
T7	Math (e.g., <i>algebra</i> , <i>geometry</i>)	Arts (e.g., <i>poetry</i> , <i>dance</i>)	Male (e.g., <i>brother</i> , <i>son</i>)	Female (e.g., <i>woman</i> , <i>sister</i>)
T8	Science (e.g., <i>experiment</i>)	Arts	Male	Female
T9	Physical condition (e.g., <i>virus</i>)	Mental condition (e.g., <i>sad</i>)	Long-term (e.g., <i>always</i>)	Short-term (e.g., <i>occasional</i>)
T10	Older names (e.g., <i>Gertrude</i>)	Younger names (e.g., <i>Michelle</i>)	Pleasant	Unpleasant

Table 1: WEAT bias tests.

From “Are We Consistently Biased? Multidimensional Analysis of Biases in Distributional Word Vectors” by Lauscher and Glavaš (2019) [\[Link\]](#)

Extrinsic Metrics – Group Fairness

- Unawareness
- Demographic Parity
- Equalized Odds
- Equality of Opportunity





Extrinsic Metrics – Equality of Opportunity

- Is only defined for binary classification
- Both classes have the same opportunity of being classified correctly
- It is measured as the difference between both classes' recalls

The background of the slide is a dense, intricate pattern of red lines representing the veins of various leaves. The lines are thin and form a complex, interconnected web across the entire frame. The colors range from a light, almost white red to a deep, dark red, creating a sense of depth and texture.

PART 3

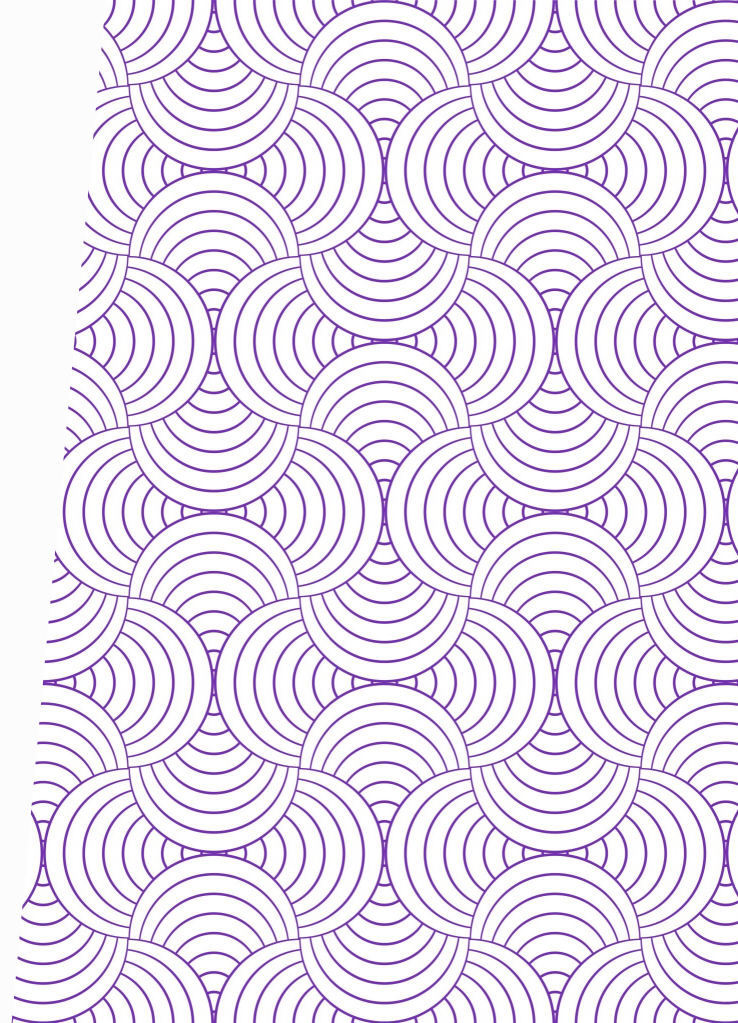
Removing Biases

Fixing Our Models to Fix the World

Removing Biases

In general there are two philosophies

- Debiasing is reducing intrinsic bias metrics
 - Note that “debiasing” is a very loaded term!
- Fairness is reducing extrinsic bias metrics
 - We will not be focusing on these for this talk

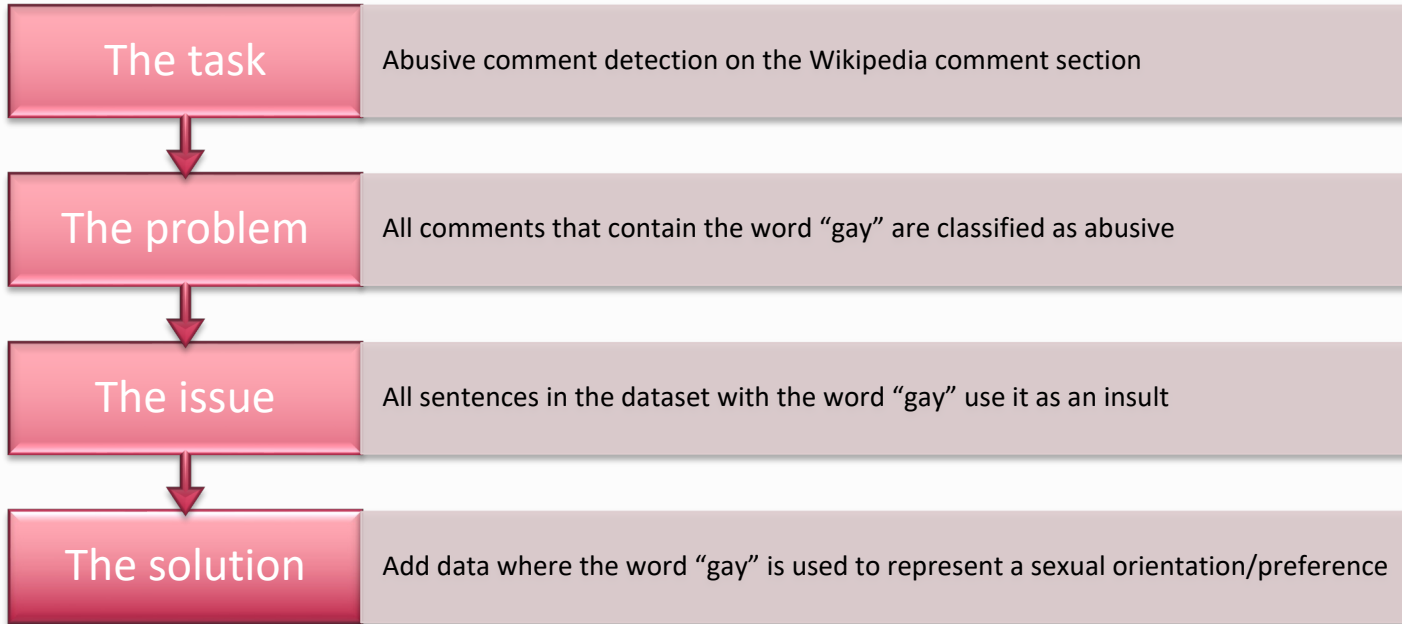


Debiasing

- **Our assumption:** removing biases in language models removes biases in downstream applications
- Can either be attempted mathematically or by altering the data



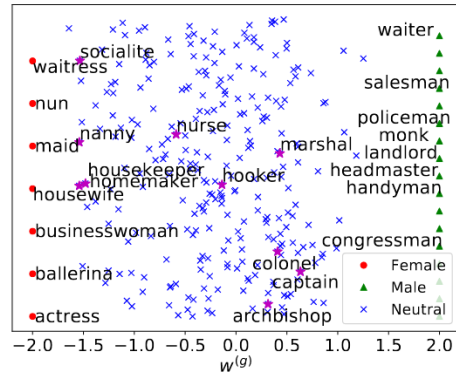
Debiasing – Dataset Balancing



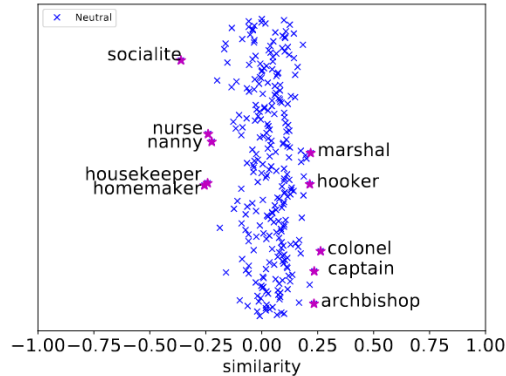
Debiasing – Removing the “Gender” Axis

- Identify the main dimension in which gender is represented in non-contextual embeddings
- Remove this dimension:
 - Completely removing it
 - Reducing the representation of non-gendered words

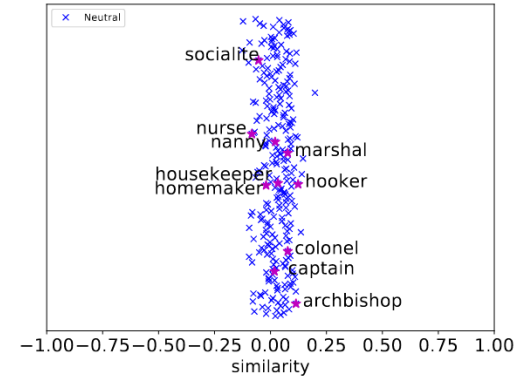
Debiasing – Removing the “Gender” Axis



(a) $w^{(g)}$ dimension for all the professions



(b) Gender-neutral profession words projected to gender direction in GloVe



(c) Gender-neutral profession words projected to gender direction in GN-GloVe

Figure 1: Cosine similarity between the gender direction and the embeddings of gender-neutral words. In each figure, negative values represent a bias towards female, otherwise male.

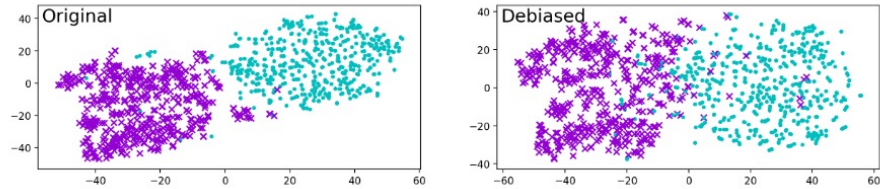
From “Learning Gender-Neutral Word Embeddings” by Zhao et al. (2018) [\[Link\]](#)

Does it Actually Work?

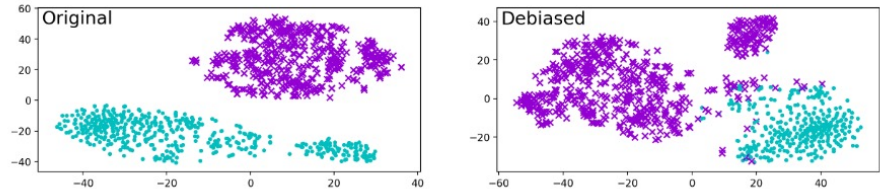
Does it Actually Work?

Not Always

Biases Might Stay Hidden



(a) Clustering for HARD-DEBIASED embedding, before (left hand-side) and after (right hand-side) debiasing.



(b) Clustering for GN-GLOVE embedding, before (left hand-side) and after (right hand-side) debiasing.

From “Lipstick on a Pig: Debiasing Methods Cover up Systematic Gender Biases in Word Embeddings But do not Remove Them” by Gonen and Goldberg (2019) [\[Link\]](#)

Biases are Complex

MODEL LEAKAGE @ F1



Predict Labels
from Image

$$p(X)$$

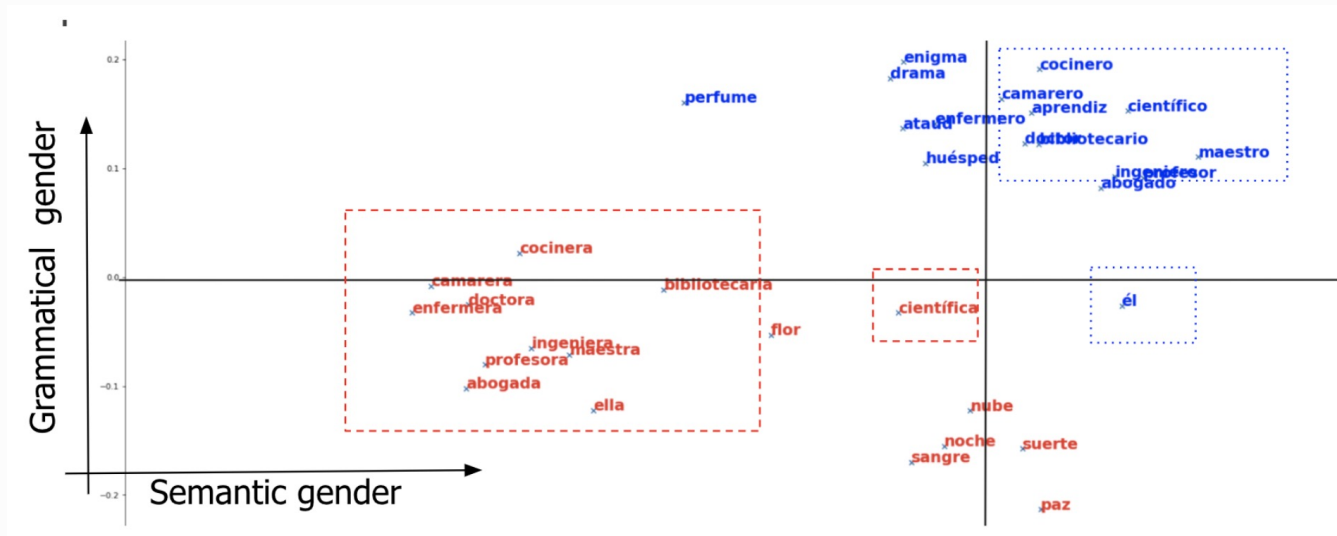
knife
spoon
fork
vase
hair drier
oven

Predict Gender from
Predicted Labels

$$f(\hat{Y})$$

From “Balanced Datasets Are Not Enough: Estimating and Mitigating Gender Bias in Deep Image Representations” by Wang et al. (2019) [\[Link\]](#)

What About Other Languages?



From “Analyzing and Mitigating Gender Bias in Languages with Grammatical Gender and Bilingual Word Embeddings” by Zhou et al. (2019) [\[Link\]](#)

A lush green forest with a stream flowing through mossy rocks. The scene is filled with vibrant green foliage, including trees and bushes, and a small stream with white water rapids. The ground is covered in moss and rocks. The overall atmosphere is serene and natural.

PART 4

But Does it Work?

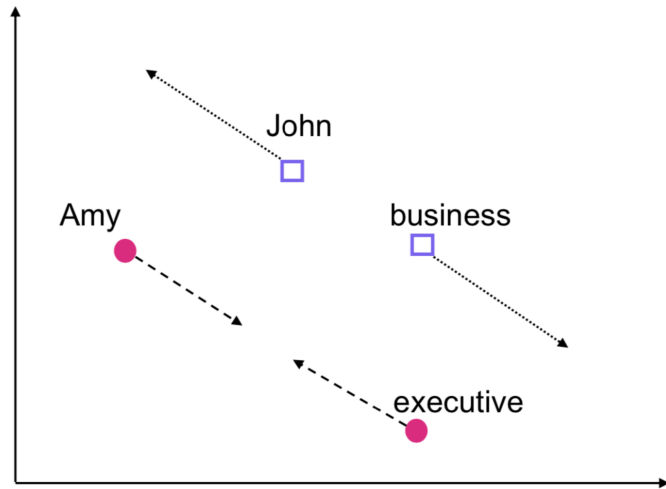
Removing Biases in Downstream Applications

General Experimental Design

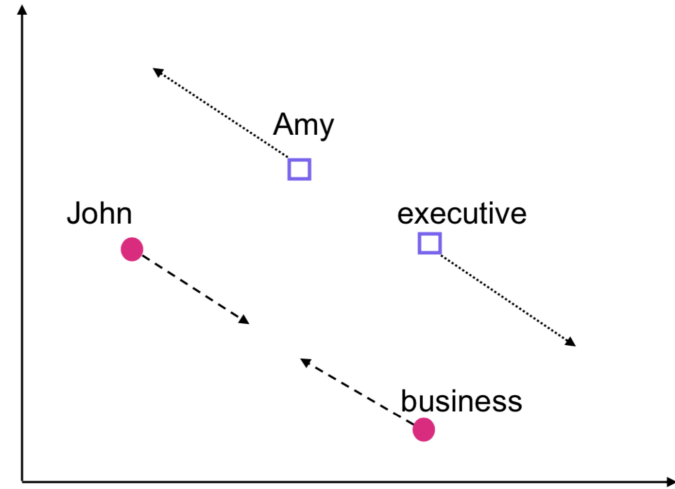
- Ways to Measure Bias
 - WEAT
 - Equality of Opportunity
- Two methods to reduce bias
 - Dataset balancing
 - Attract-Repel
- Two and a half downstream applications
 - Correference resolution in English
 - Hatespeech detection in English and in Spanish



Attract-Repel



Debias



Overbias

About the Languages

English

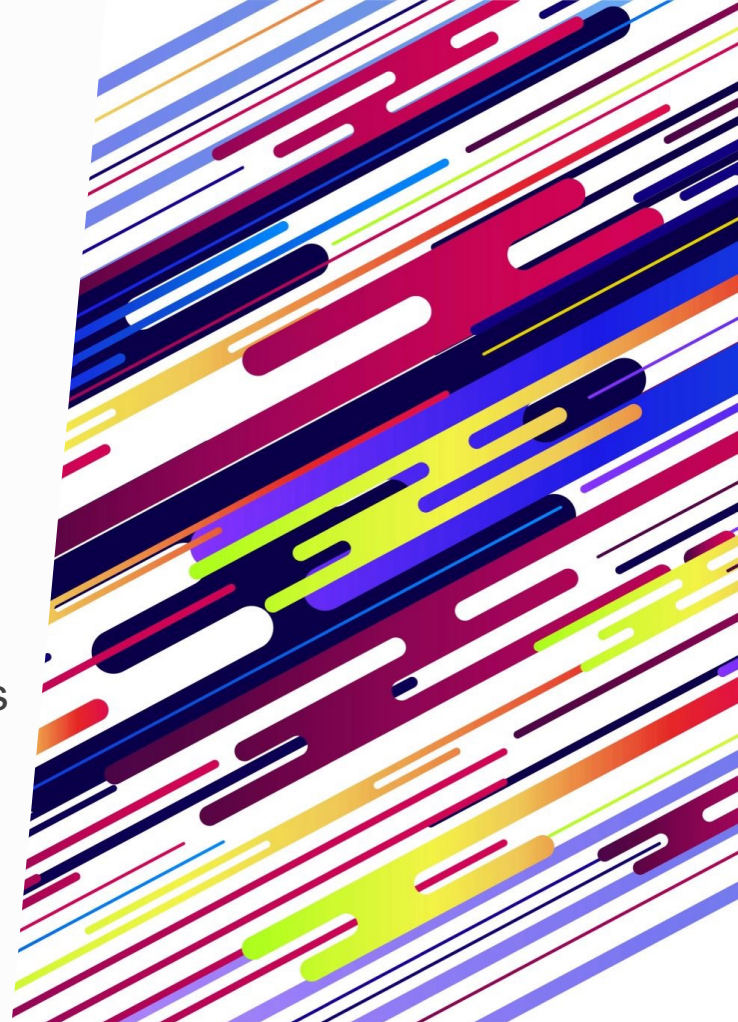
- Has been used in most bias studies
- Only has semantic gender

Spanish

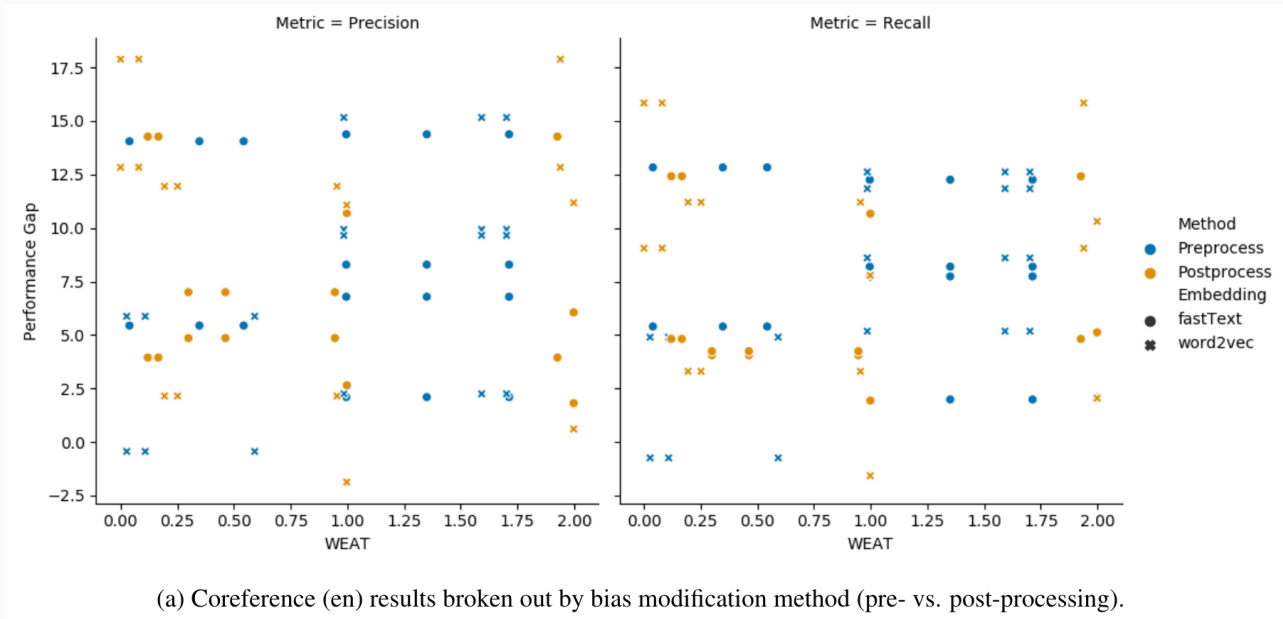
- Some studies have analysed biases in Spanish embedding spaces
- There is enough data available for most tasks
- Has both grammatical and semantic gender

Downstream Applications

- Correference resolution for gendered pronouns
 - A stereotypical task for bias assessment in English
 - However, it's trivial in Spanish
- Hate speech classification
 - Allows us to compare the effects of semantic vs grammatical bias
 - Against women (English and Spanish)
 - Against migrants (Spanish)

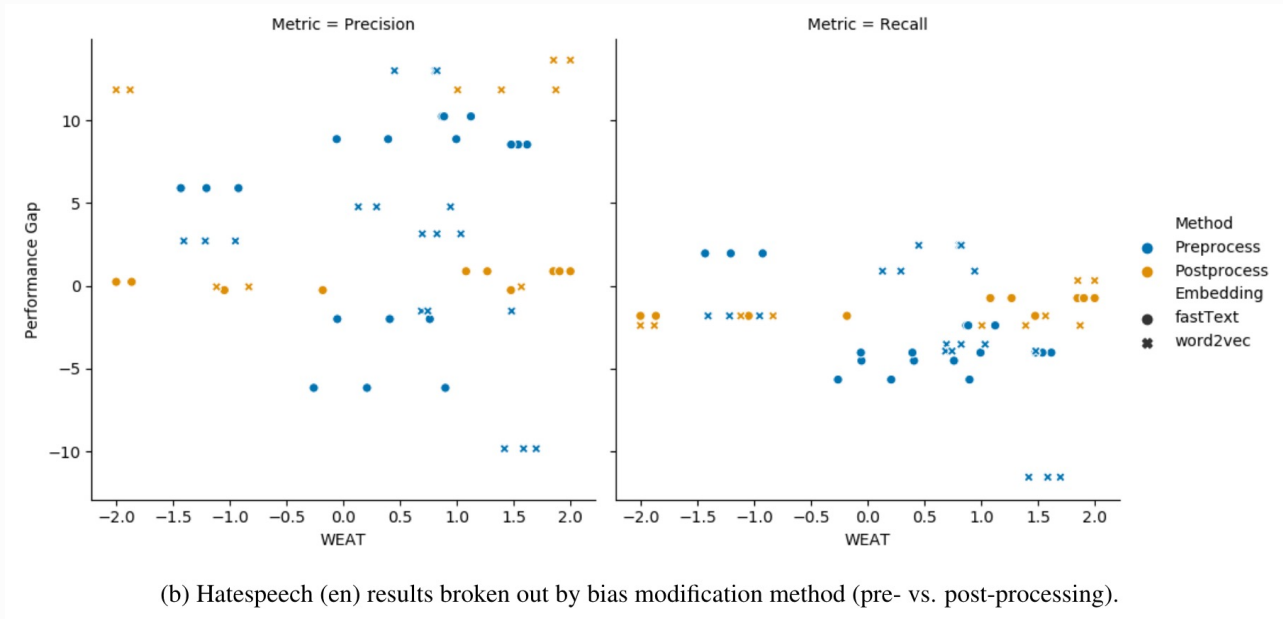


The Results – Correferrence (eng)



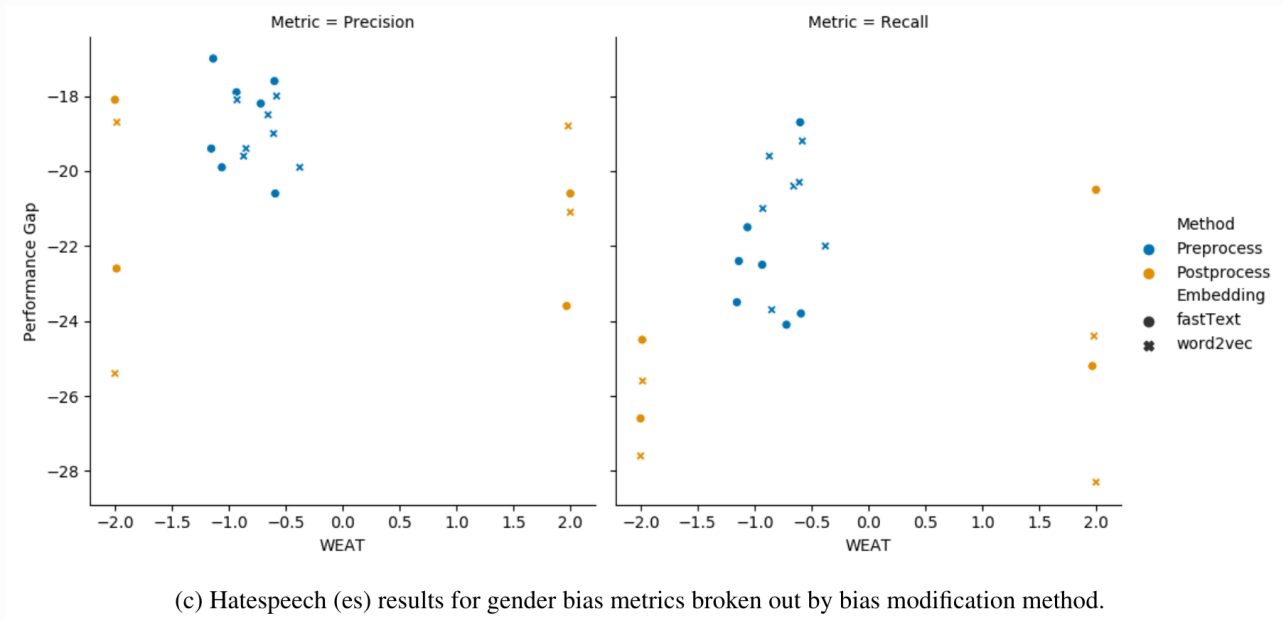
From “Intrinsic Bias Metrics Do Not Correlate with Application Bias” by Goldfarb-Tarrant et al. (2021) [\[Link\]](#)

The Results – Hate Speech Detection (eng)



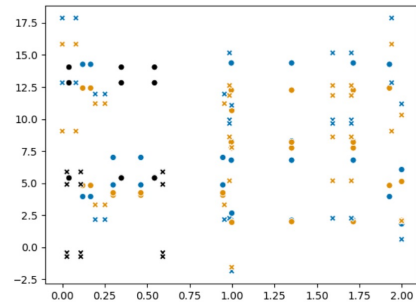
From “Intrinsic Bias Metrics Do Not Correlate with Application Bias” by Goldfarb-Tarrant et al. (2021) [\[Link\]](#)

The Results – Hate Speech Detection (spa)

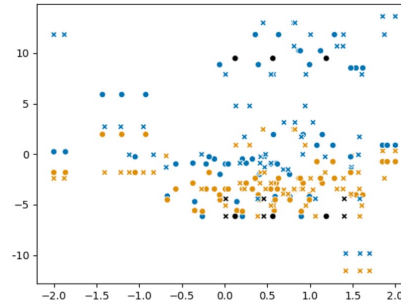


From “Intrinsic Bias Metrics Do Not Correlate with Application Bias” by Goldfarb-Tarrant et al. (2021) [\[Link\]](#)

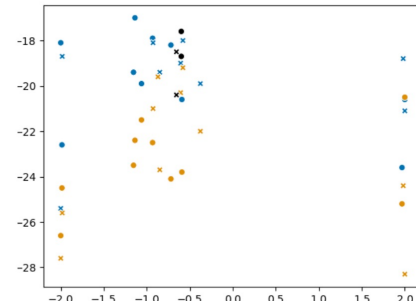
The Results



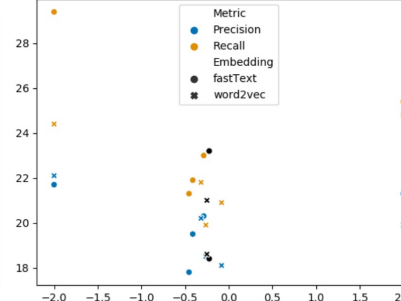
(a) Coreference results (en), gender bias



(b) Hatespeech detection results (en), gender bias



(c) Hatespeech detection results (es), gender bias



(d) Hatespeech detection results (es), migrant bias

From “Intrinsic Bias Metrics Do Not Correlate with Application Bias” by Goldfarb-Tarrant et al. (2021) [\[Link\]](#)



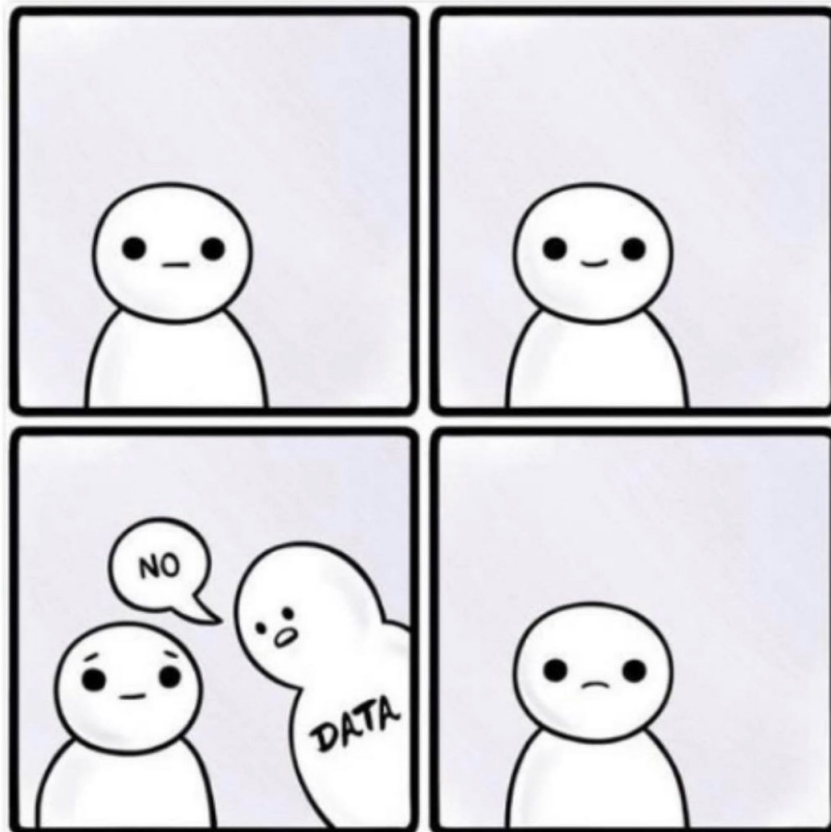
Dorsa Amir

@DorsaAmir

Follow

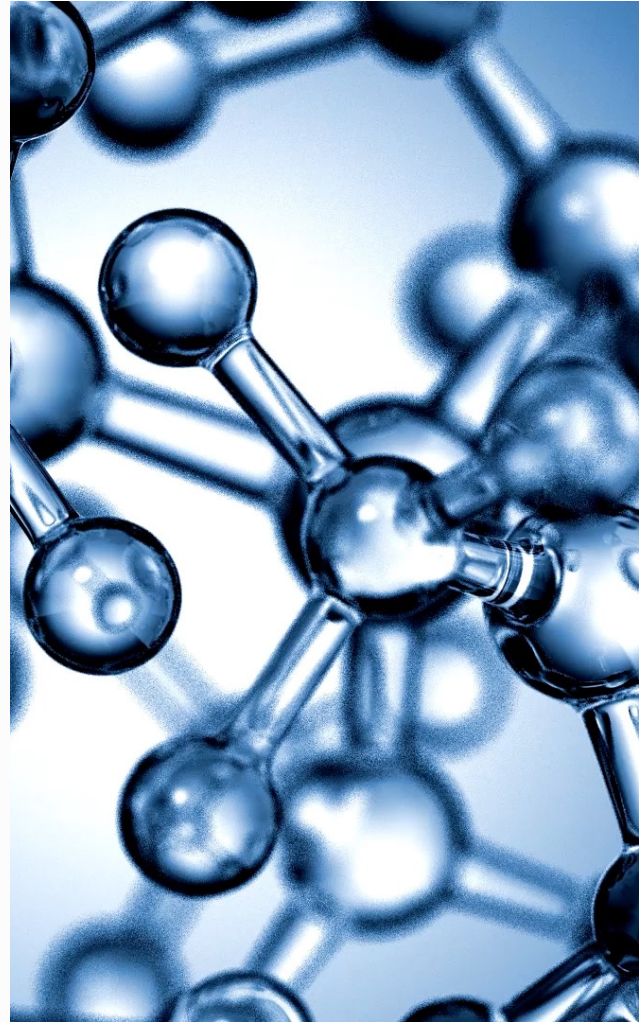


The (real) scientific method.



Our Insights

- Reducing bias in embedding spaces is unpredictable in terms of downstream application biases
 - It has been reproduced with more downstream applications
 - Language models also have these issues
- Spanish (X)WEAT is biased itself!
 - Almost all science words were grammatically male
 - Some issues with translations
 - No usual names from Spanish-speaking countries



Going Forward

- Most bias and fairness research focuses on
 - Gender as a binary (male/female)
 - Race in the United States as a binary (white/black)
- Biases are very diverse but the experiments ran more often than not aren't
- Getting this research into the hands of those who need it is important





GÖTEBORGS
UNIVERSITET

SPRÅKBANKEN TEXT

Ricardo Muñoz Sánchez

ricardo.munoz.sanchez@svenska.gu.se

[rimusa.github.io](https://github.com/rimusa)