

Hidden Biases

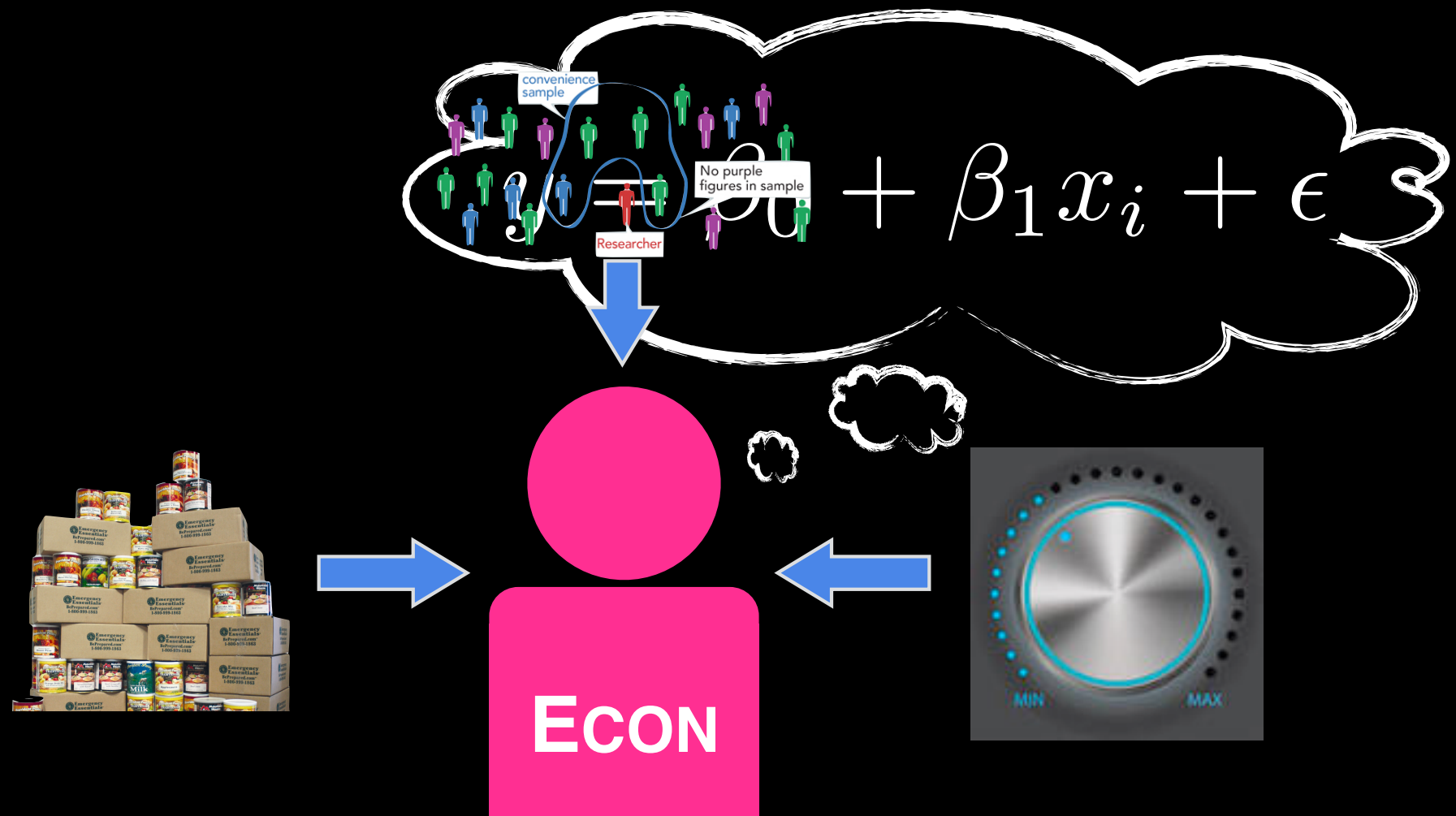
Ethical Issues in NLP, and What to Do about Them

Dirk Hovy

dirk.hovy@unibocconi.it

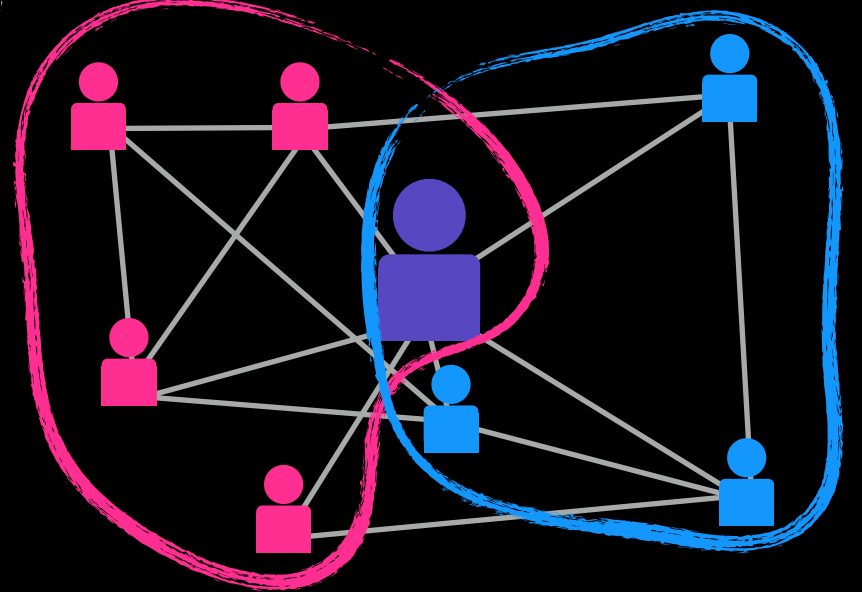
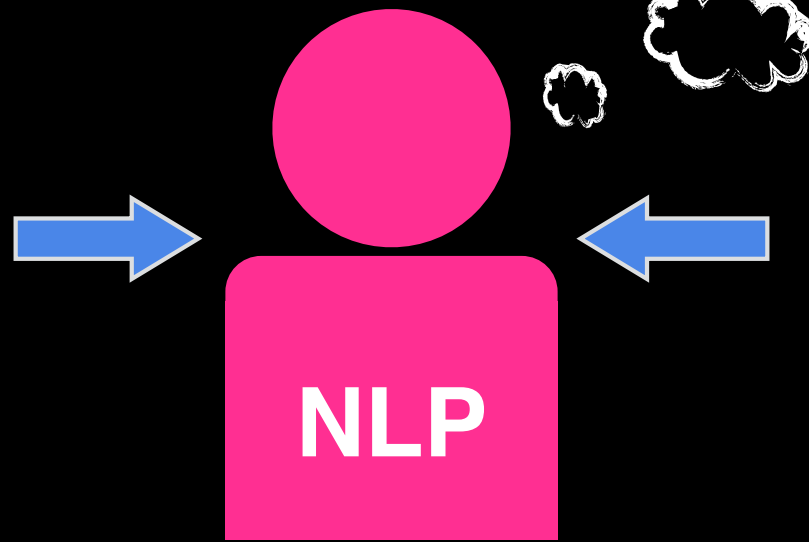
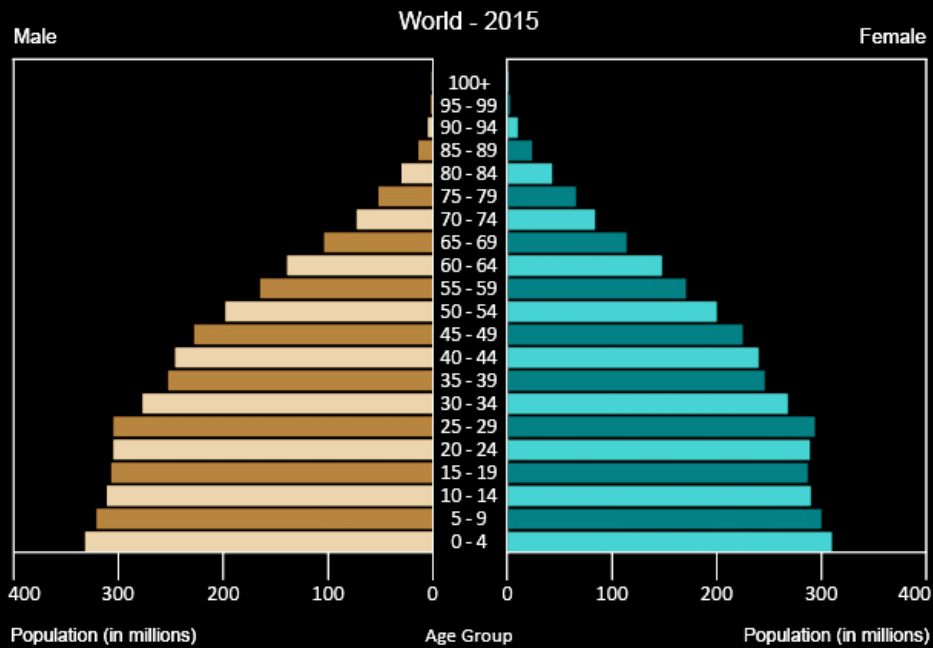
 @dirk_hovy

A Limited View

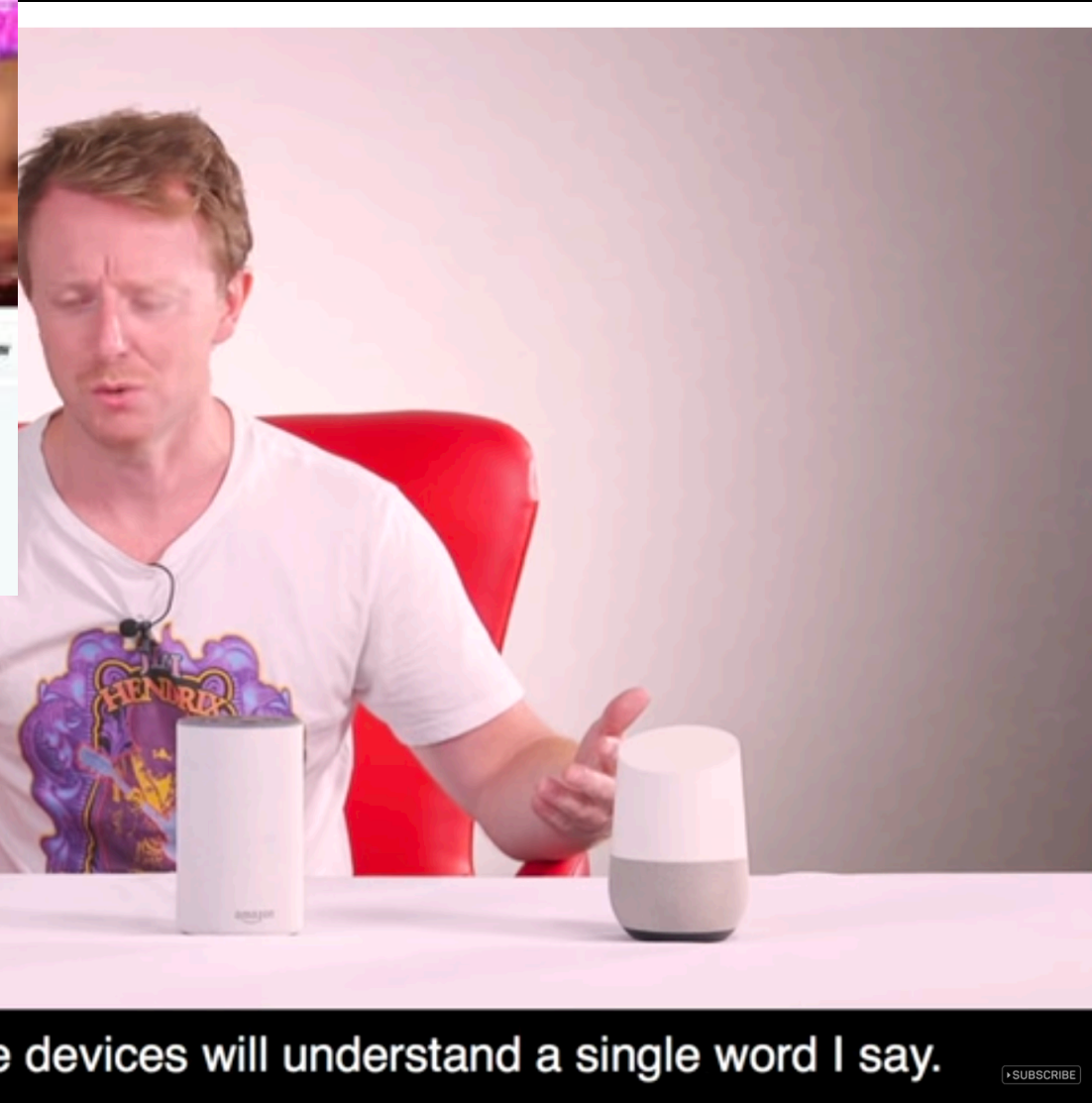


Language as Information

$$y = \beta_0 + \beta_1 x_i + \epsilon$$



Biased Language Systems

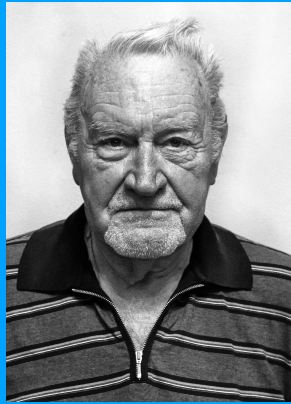


It's shite being Scottish in a smart speaker world

70,140 views

1.7K 118 SHARE SAVE ...

Language Biases



Example 1

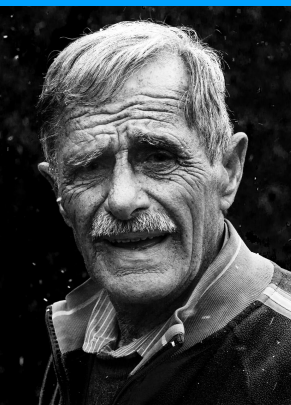
I don't understand you...



Example 2

System

Hello,
computer

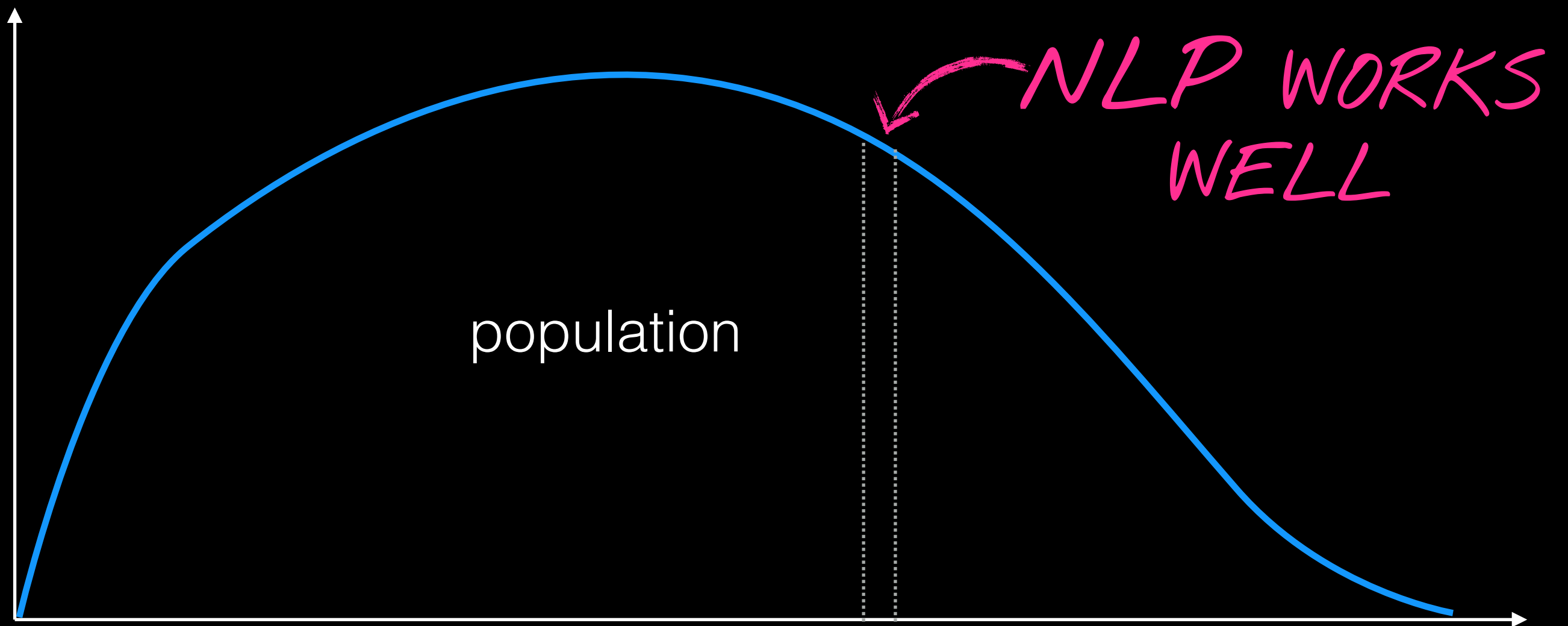


Example N

Shite...



The Consequences



Solutions?

SYMPTOMS



Example 1

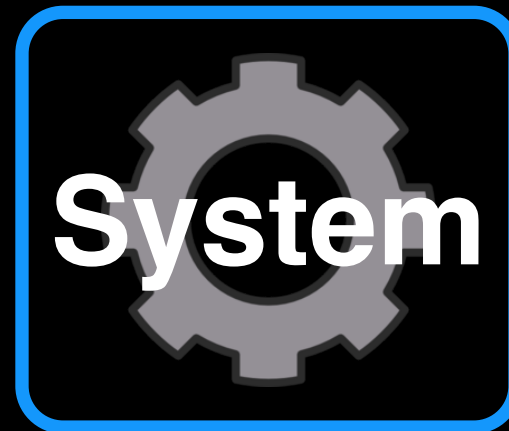


Example 2



Example N

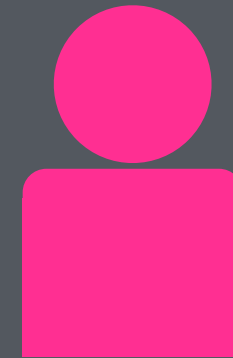
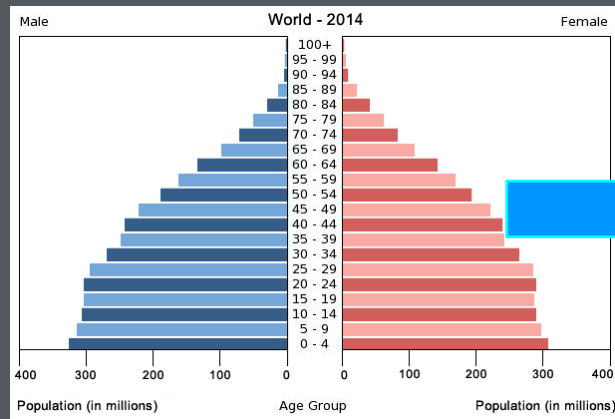
CAUSE



Goals for Today

- Point out potential **ethical issues** in NLP
- Introduce 4 **sources of bias**
- Discuss **counter measures**

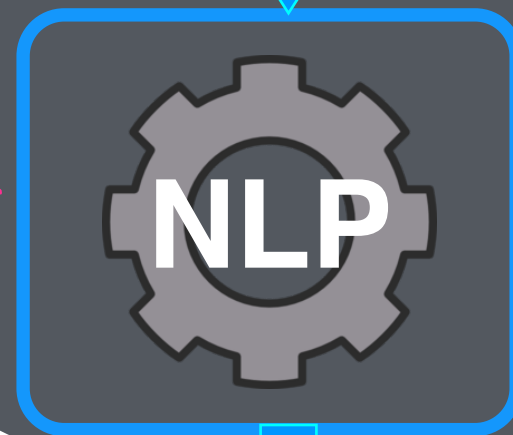
Sources of Bias



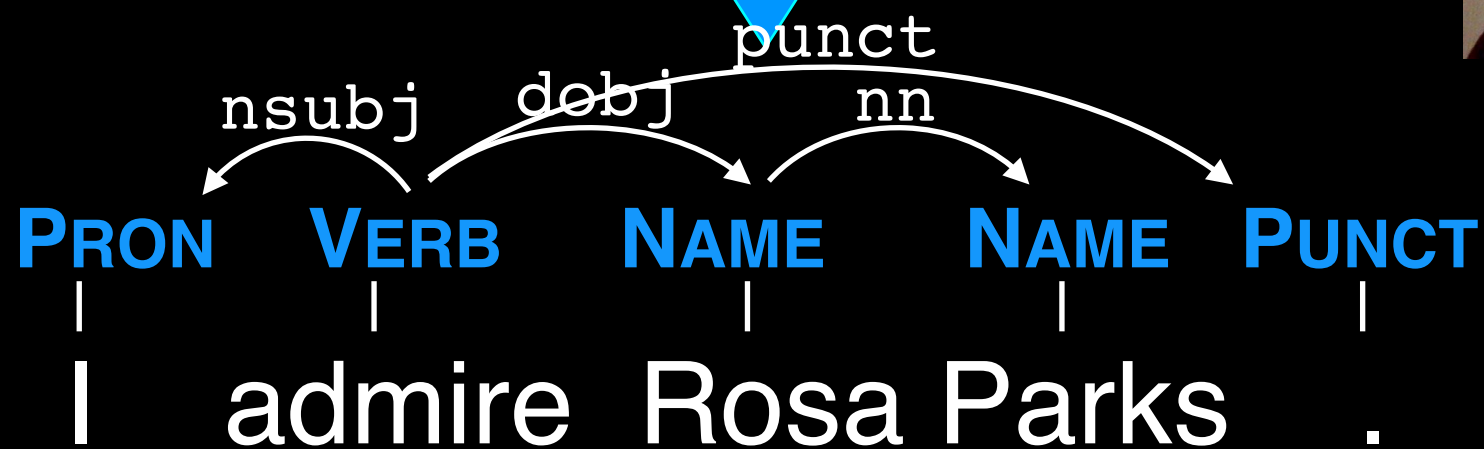
SELECTION

ANNOTATION

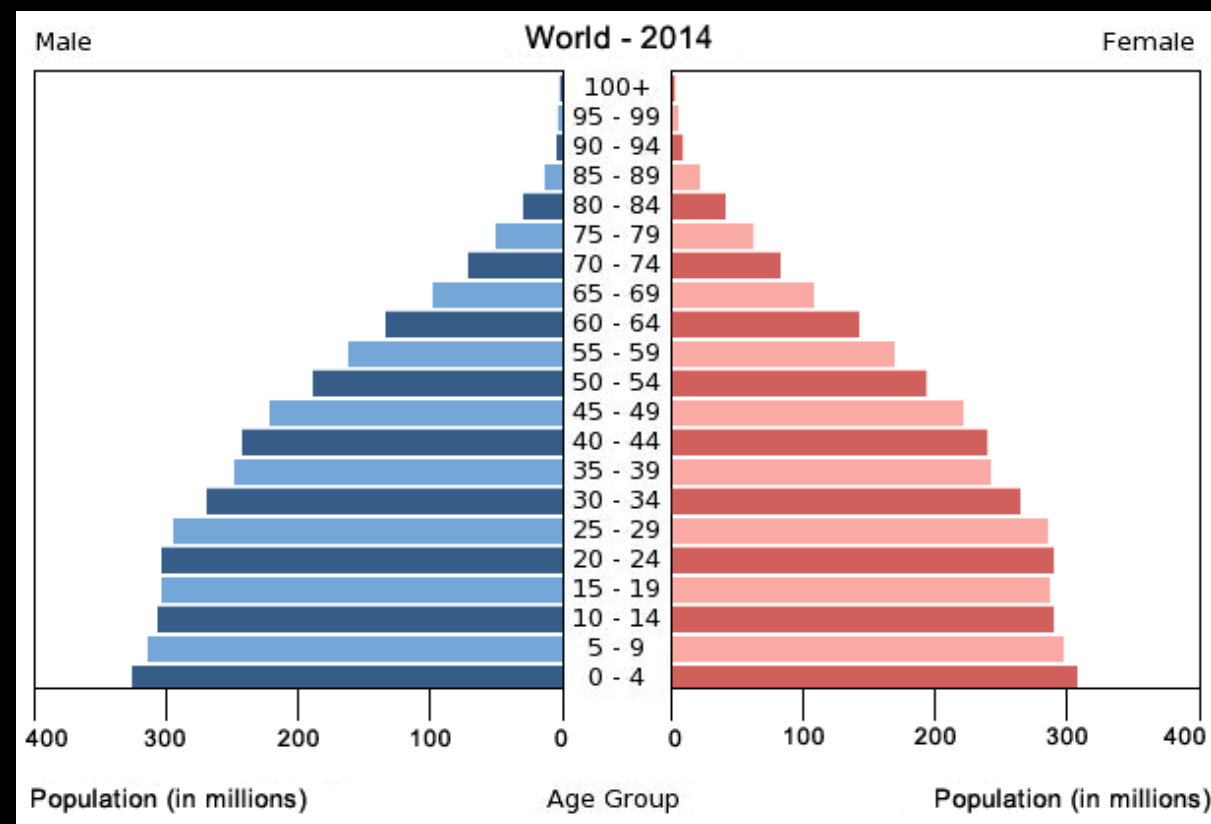
MODELS



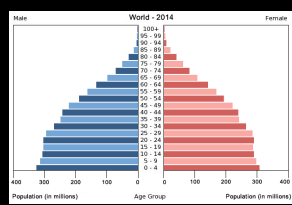
DESIGN



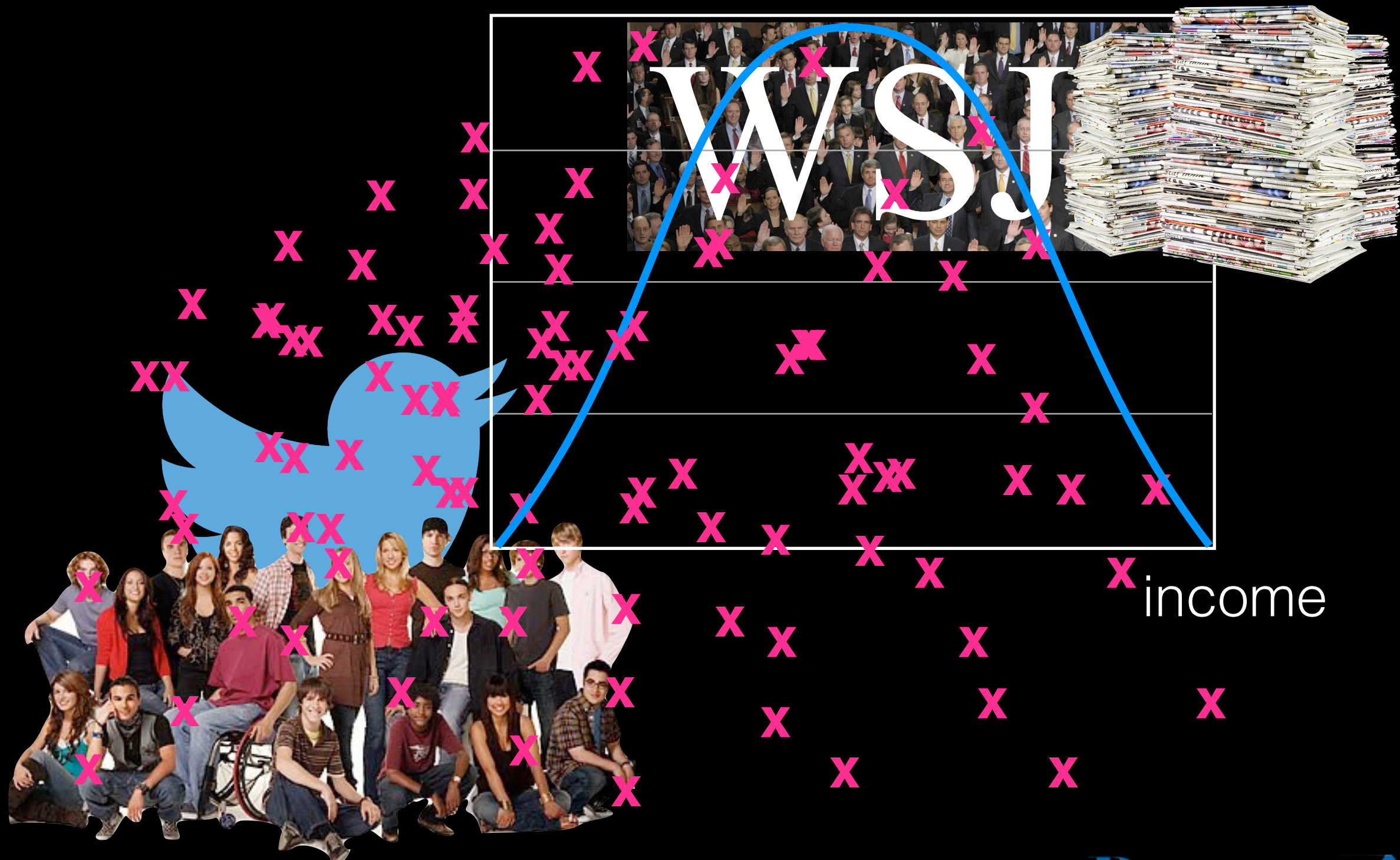
Part 1: Data Bias



Distributions

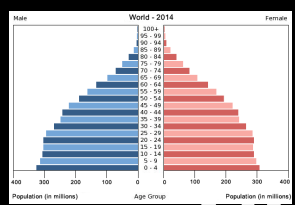


age



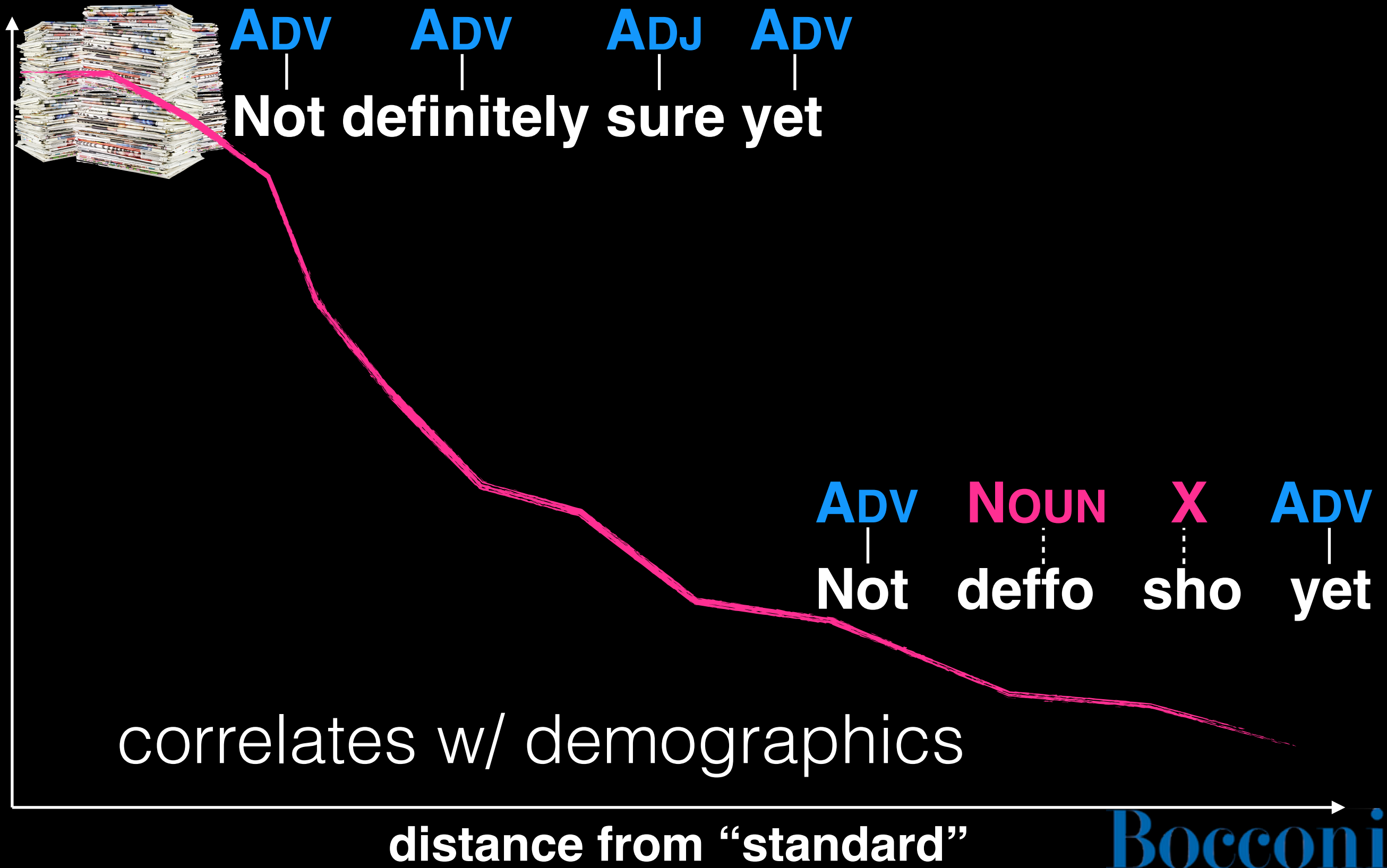
income

The WSJ Effect

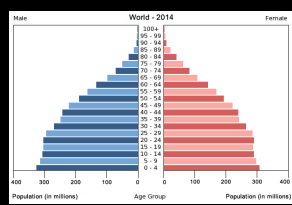


NLP

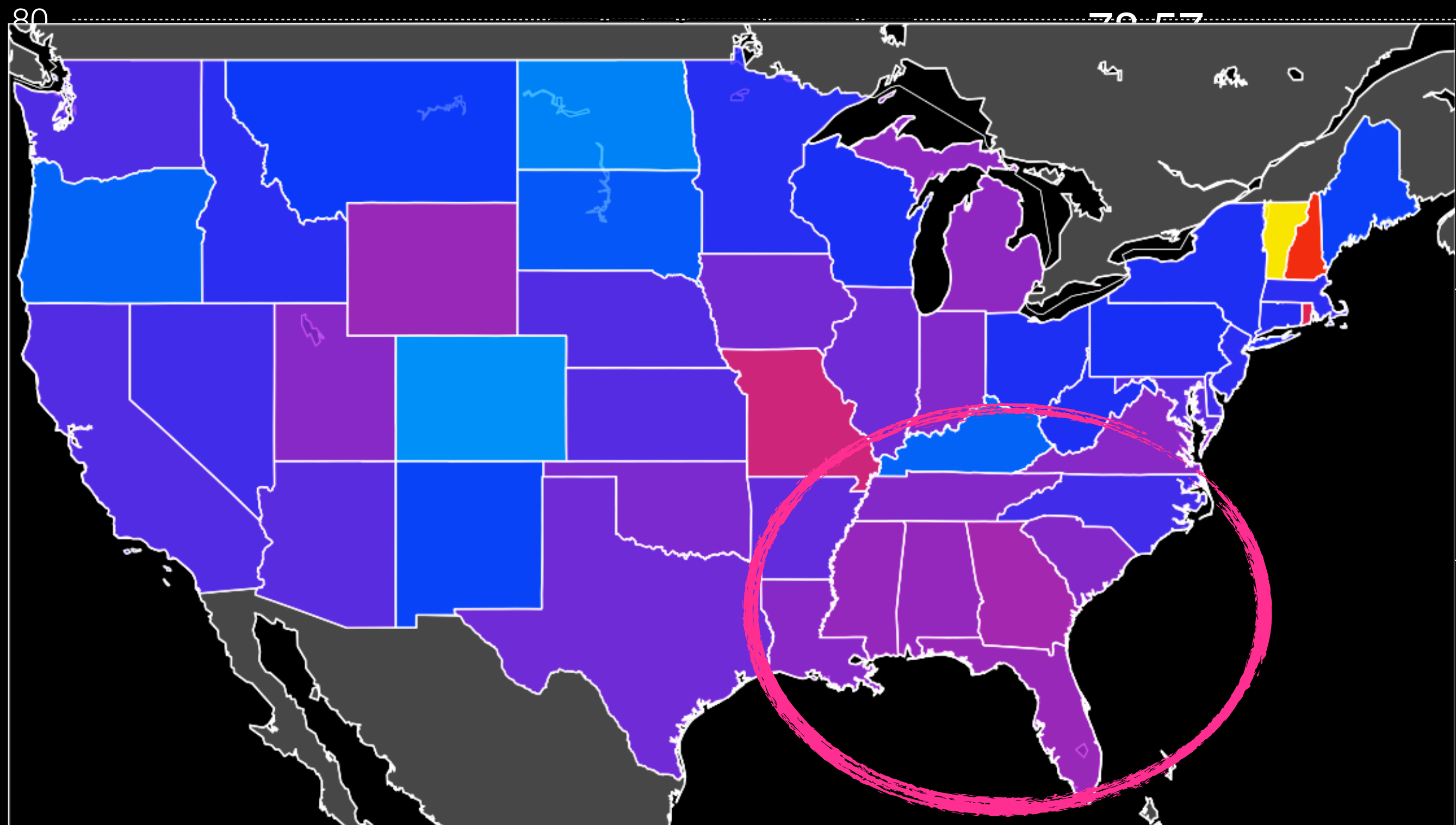
performance



Exclusion



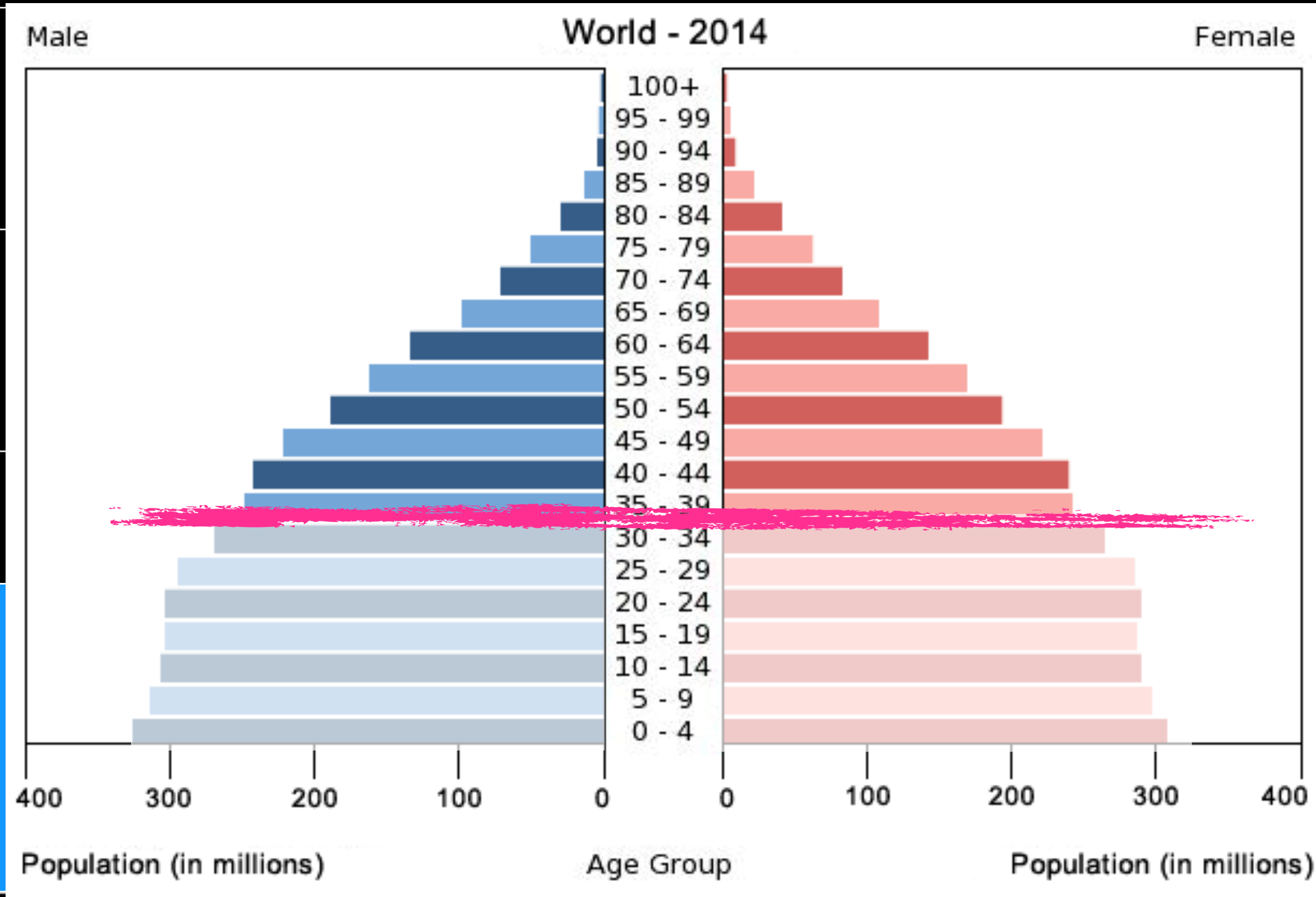
F1



Exclusion

accuracy

100



ing
views
ned

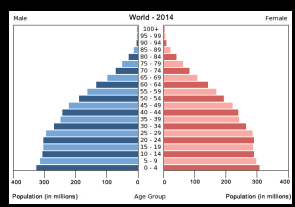
2

O45

U35

O45

U35

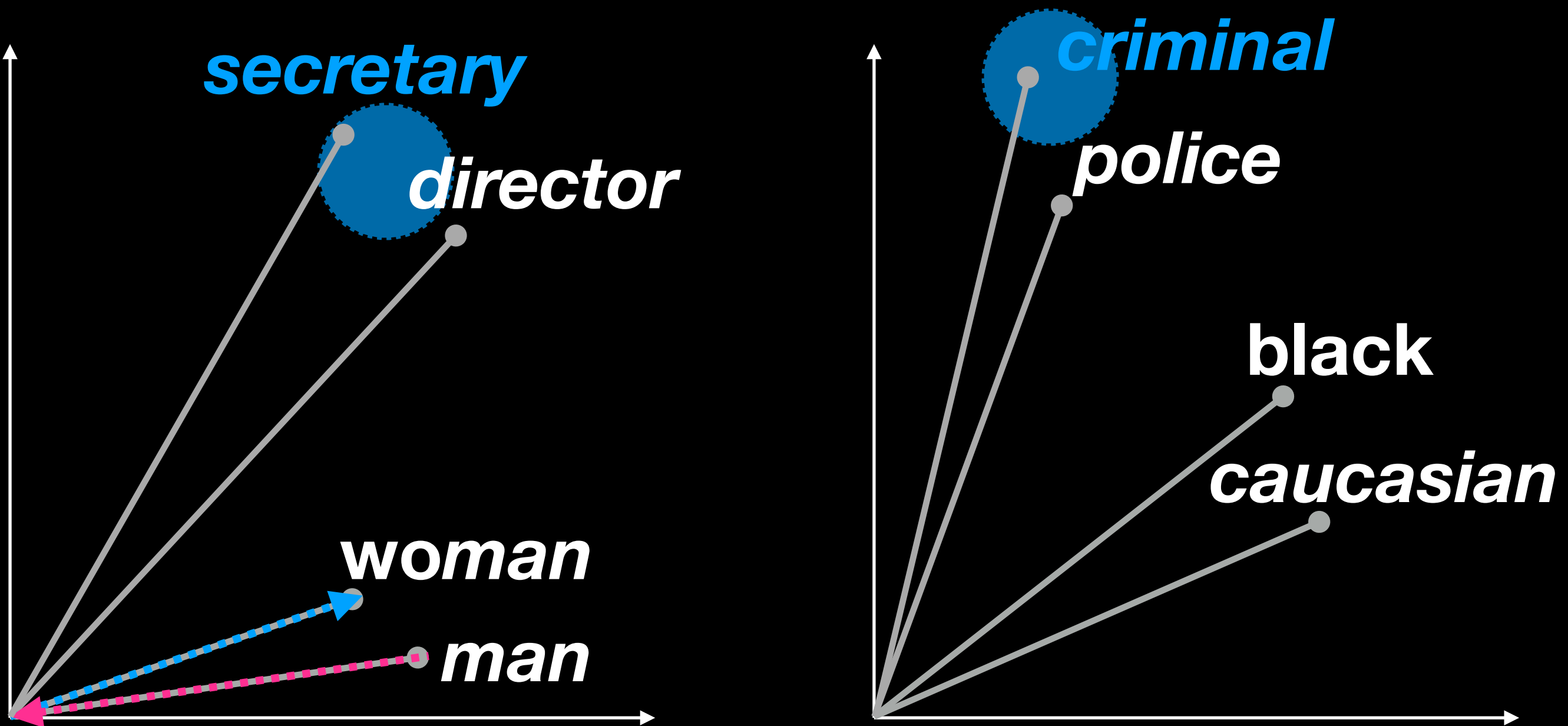


Biased Vectors

Bolukbasi et al. (2016)

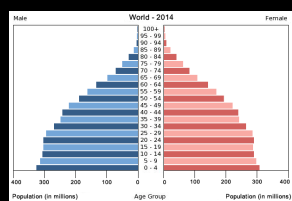
Manzini et al. (2019)

director – *man* + *woman* \approx ***secretary***
police – *caucasian* + *black* \approx ***criminal***

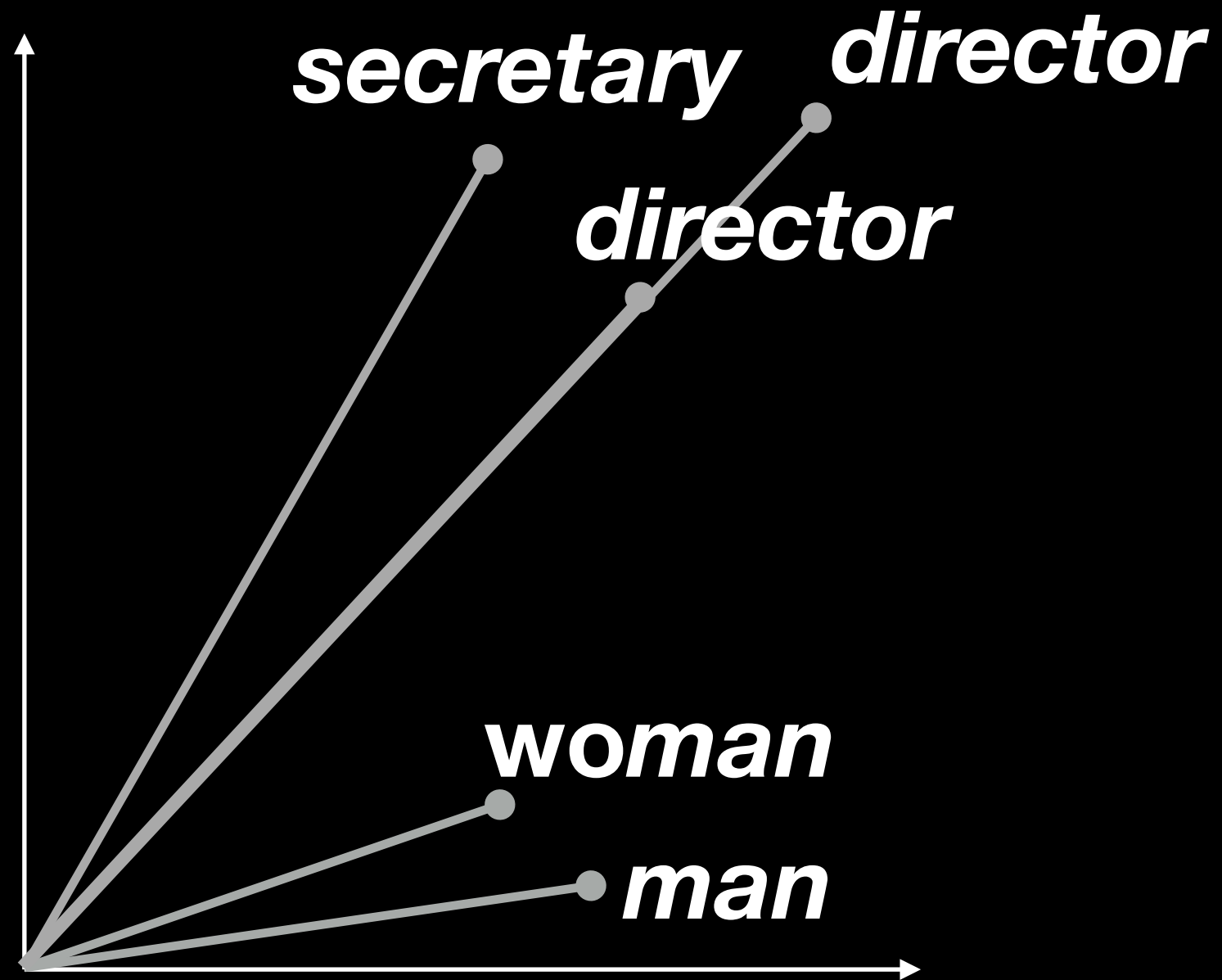
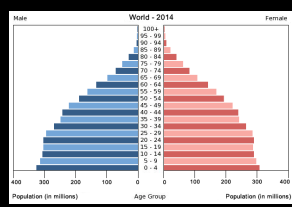


Idea!

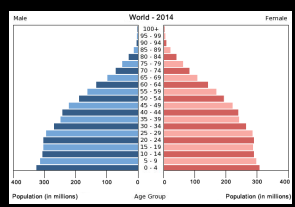
DEBIAS THE VECTORS!



Debiasing Vectors

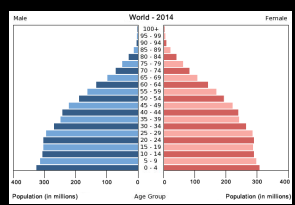


Cause vs. Symptoms

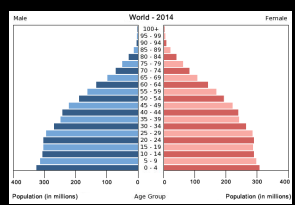


Idea!

INCLUDE DEMOGRAPHIC INFORMATION
IN TEXT REPRESENTATION



Systems



AGNOSTIC

INFORMED

training data

This is a tiny little example text written by someone.

This is a tiny little example text written by someone.

training data

This is a tiny little example text written by someone.

This is a tiny little example text written by someone.

training data

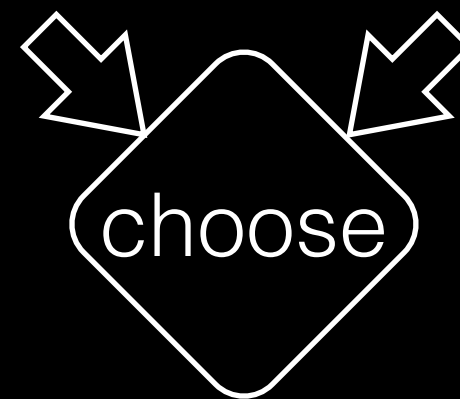
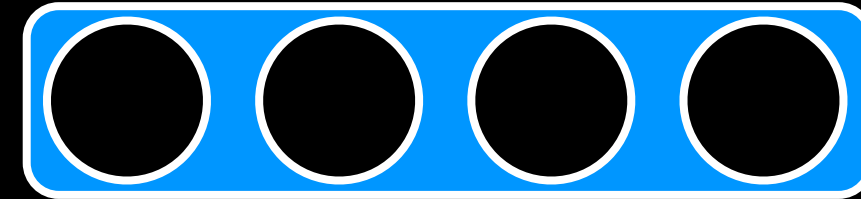
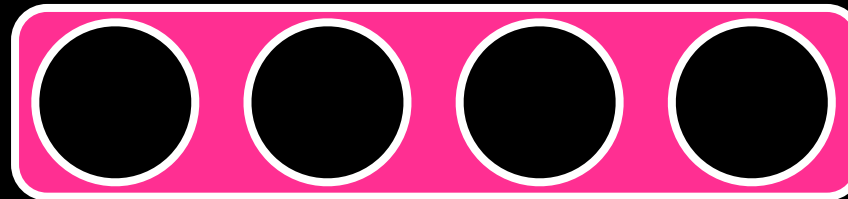
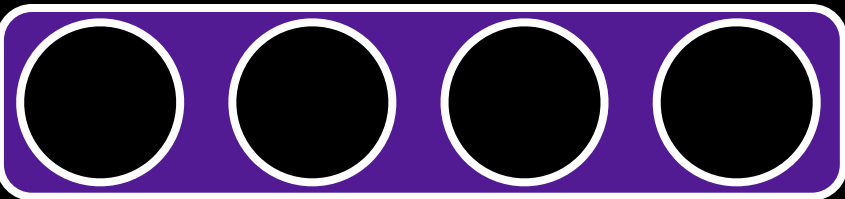
This is a tiny little example text written by someone.

This is a tiny little example text written by someone.

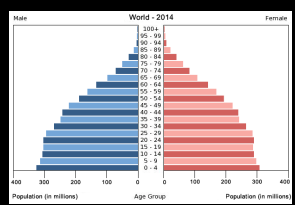
+

+

+



Results for Age (avg)



F1

65

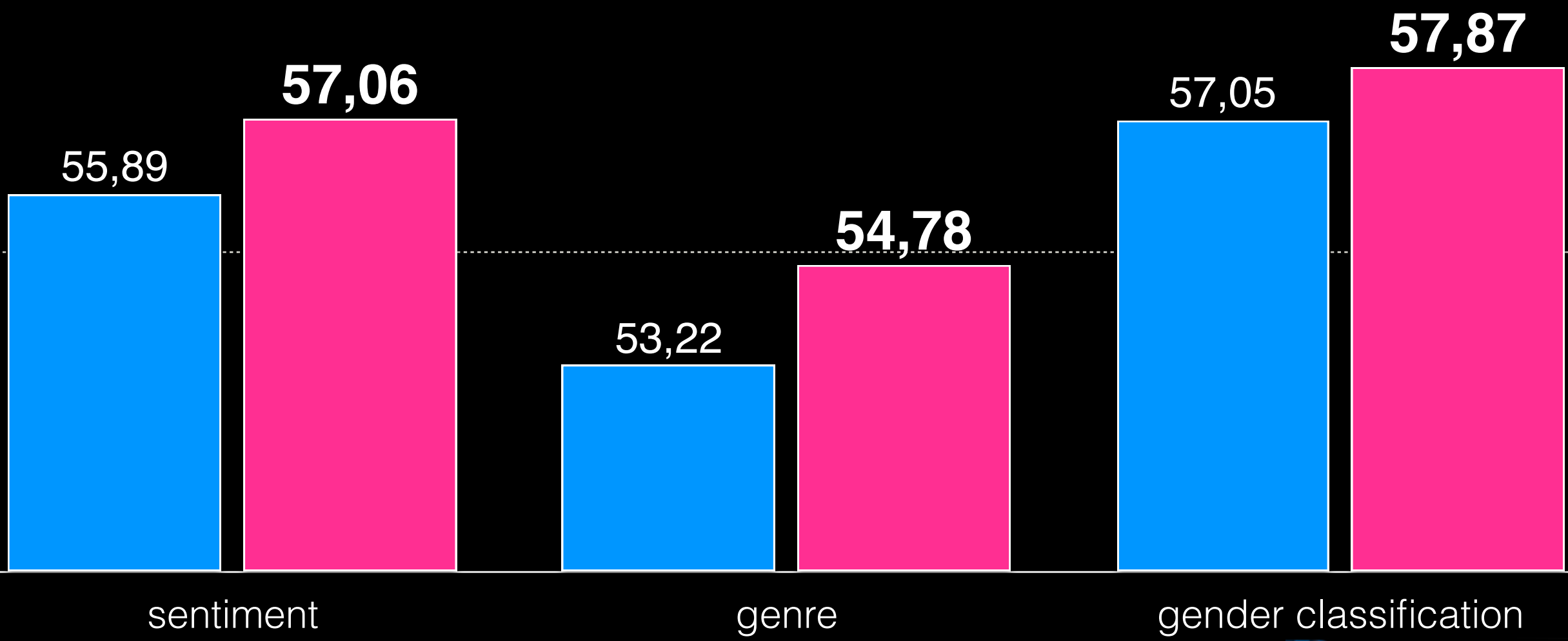
agnostic
aware

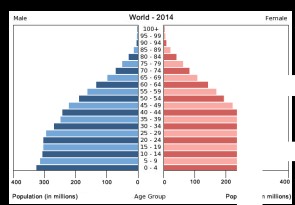
60

55

50

21





Results for Gender (avg)

F1

65

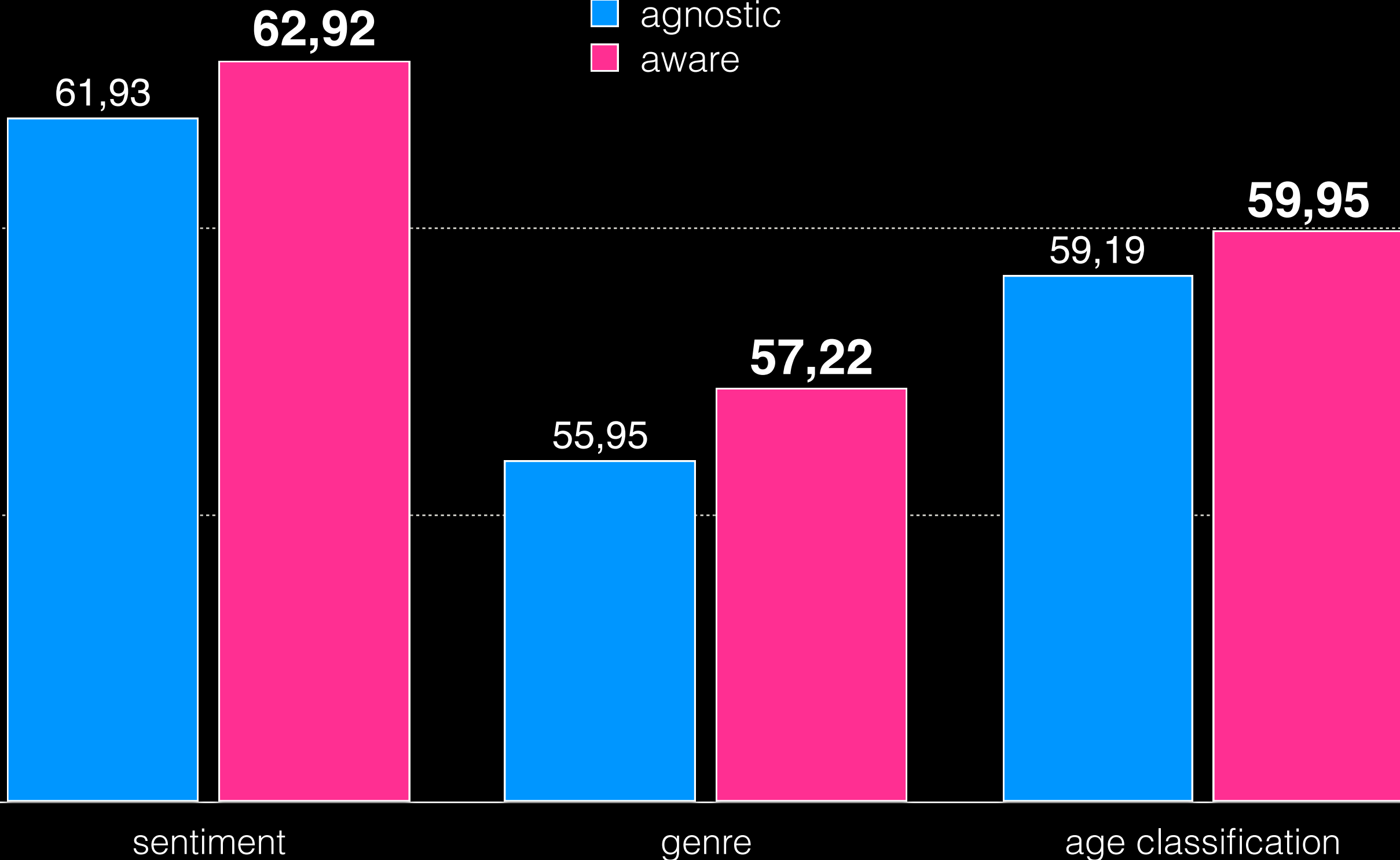
60

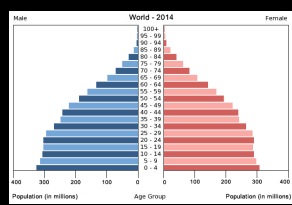
55

50

22

agnostic
aware



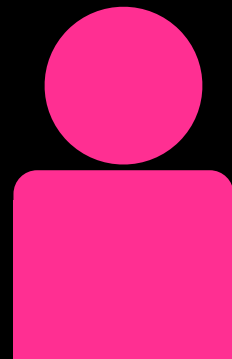


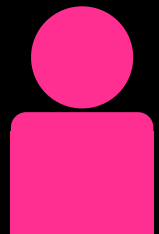
Ok, but...



WHAT IF WE DON'T KNOW
WE WANT TO KNOW THE
AUTHOR'S DEMOGRAPHICS?

Part 2: Annotation Bias





Annotator Bias



It's a
particle!



No! It's an
adposition!

PRON VERB

PRT

NOUN NUM

PRON VERB

ADP

NOUN NUM

it comes out apr 30

Idea!

FIND OUT WHO'S RELIABLE!



Model

TRUTH

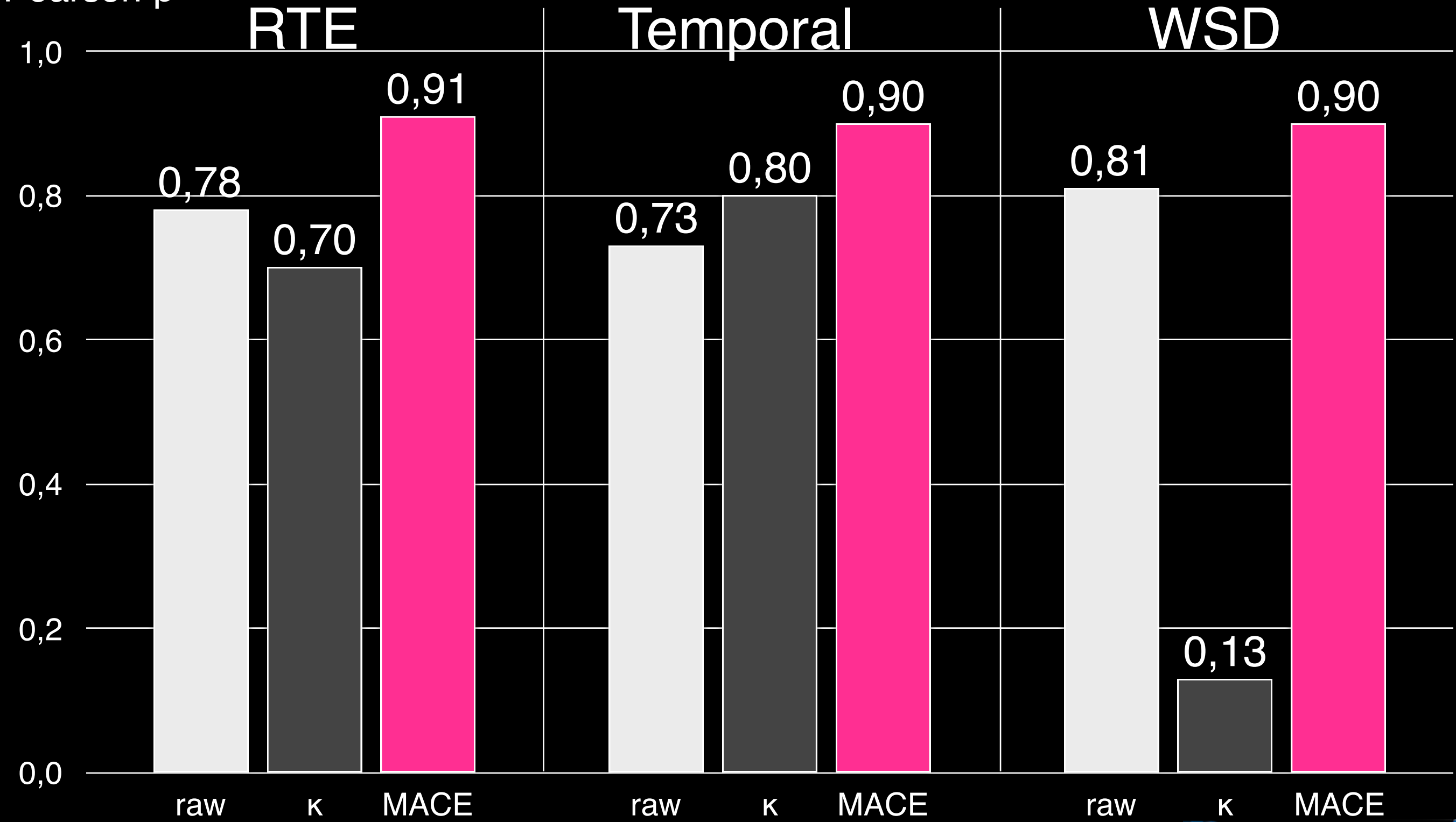


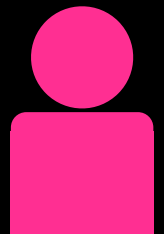
www.isi.edu/publications/licensed-sw/mace/

www.dirkhovy.com/portfolio/papers/download/mace.zip

Correlation with Proficiency

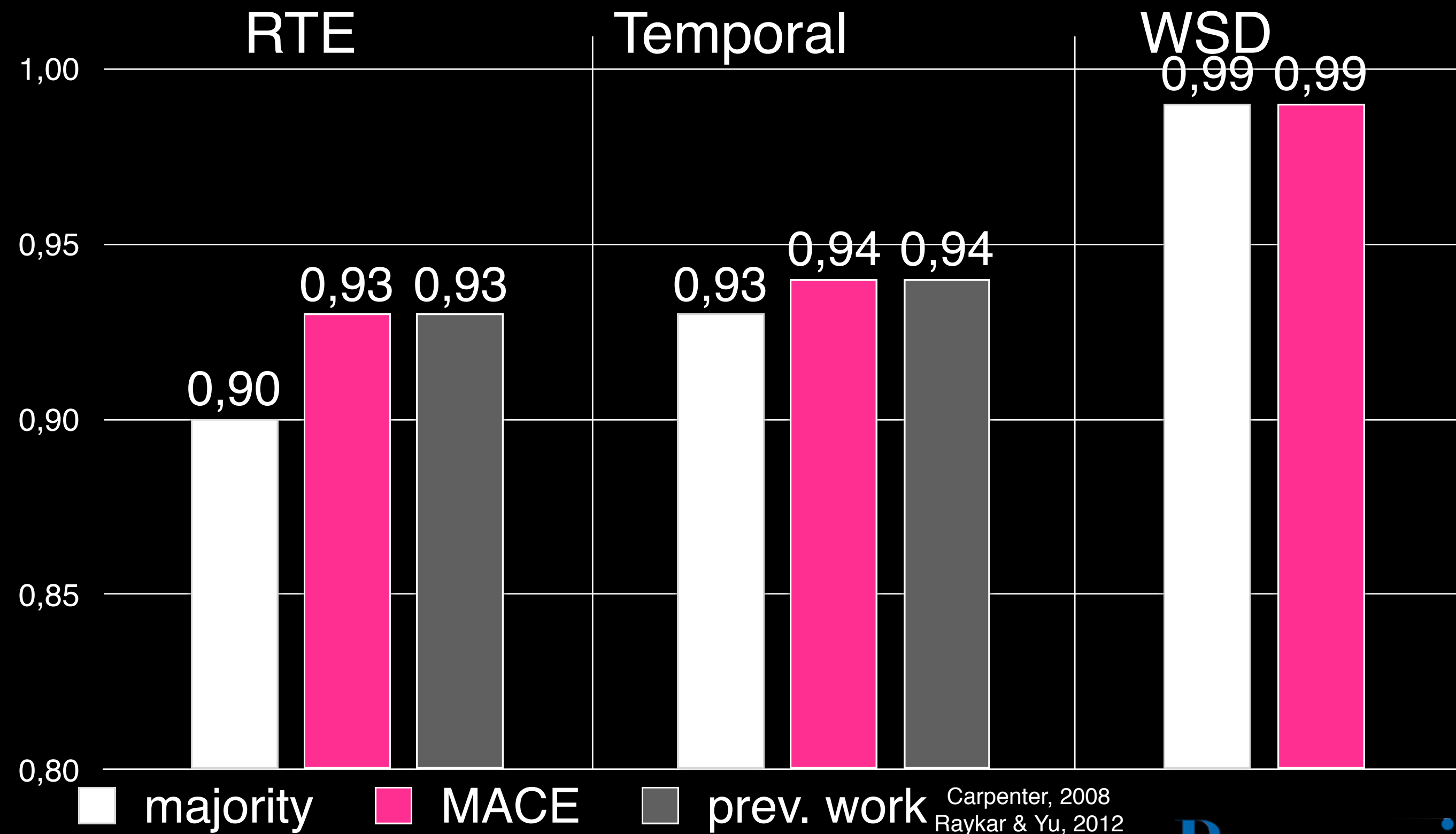
Pearson ρ





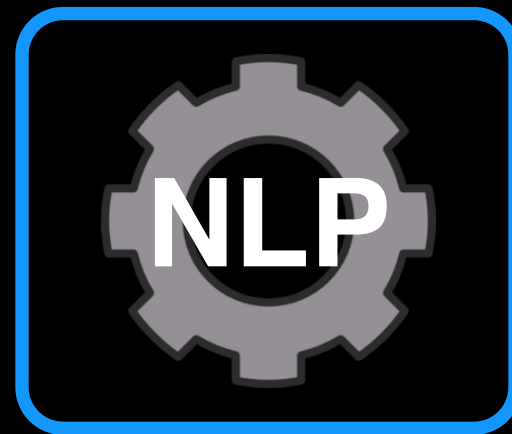
Prediction Accuracy

accuracy



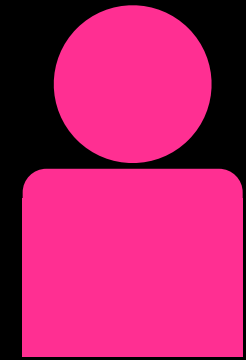
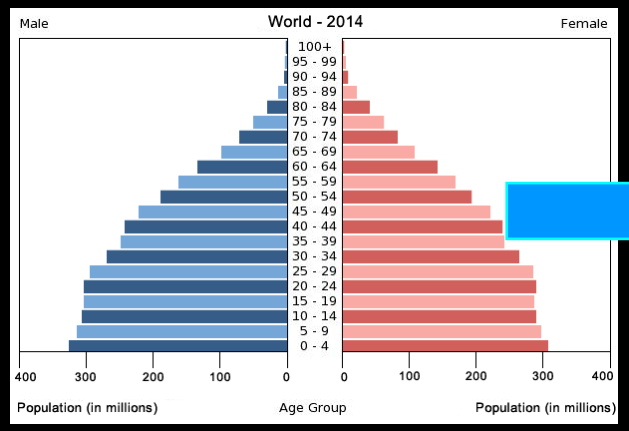
majority MACE prev. work Carpenter, 2008
Raykar & Yu, 2012

Part 3: Model Bias



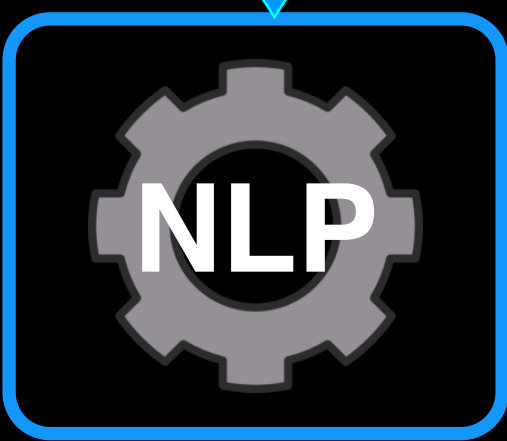


Biased Models



SELECTION

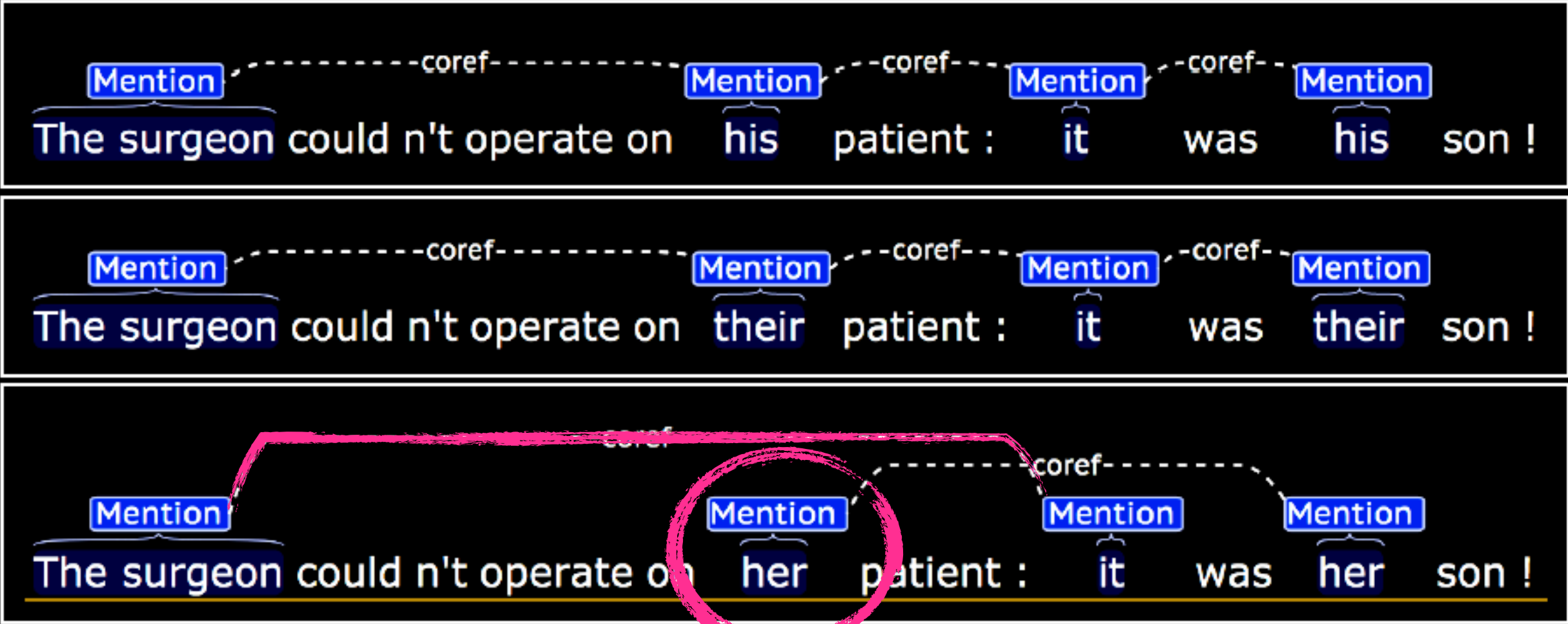
ANNOTATION



THIS IS REPRESENTATIVE

THIS IS RELIABLE

Wrong Coreference



Biased Sentiment Analysis

0.64

He made me feel **afraid**

0.52

I made **Latisha** feel **angry**

0.48

She made me feel **afraid**

0.43

I made **Heather** feel **angry**



Models Amplifying Bias

BIAS = 0.66



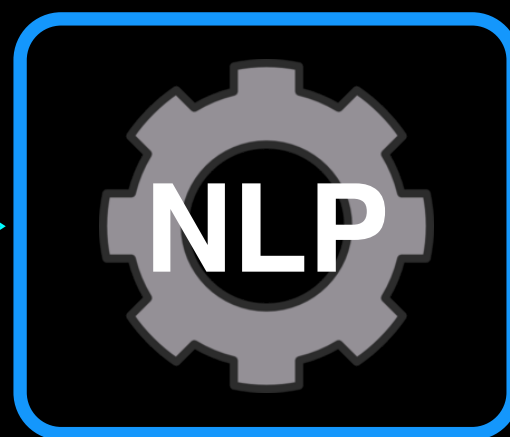
Agent: WOMAN



Agent: MAN



Agent: WOMAN



BIAS = 0.84



Agent: WOMAN



Agent: WOMAN



Agent: WOMAN



Agent: MAN



Agent: WOMAN



Idea!

*DISCOURAGE MODELS FROM
AMPLIFICATION!*





Reducing Bias

BIAS = 0.66



Agent: WOMAN

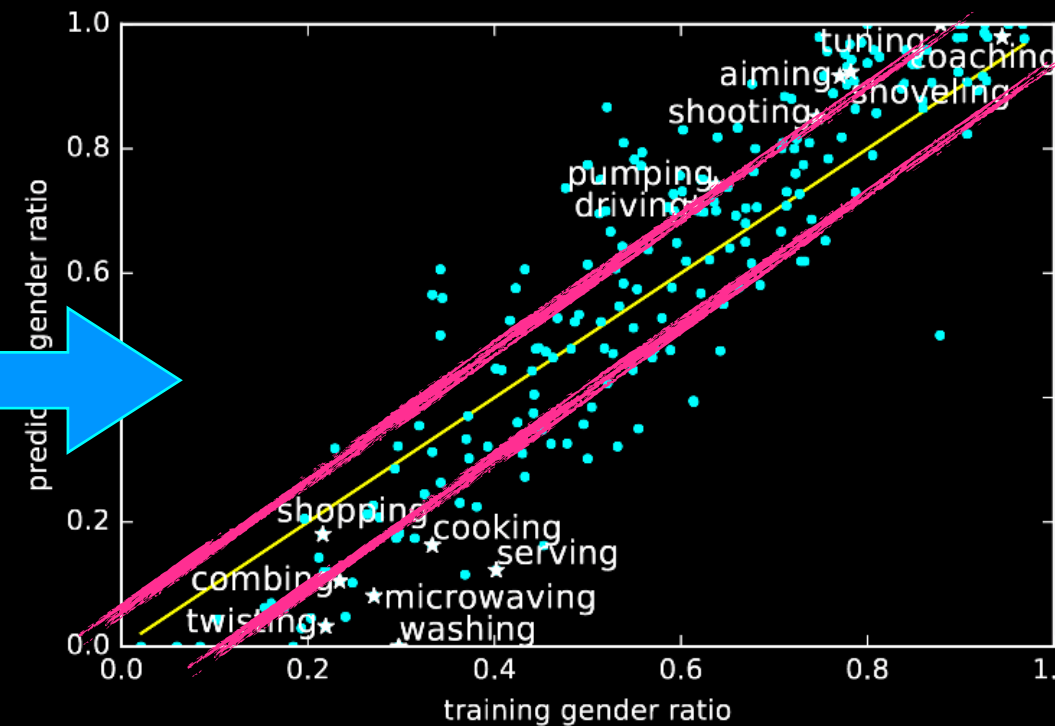


Agent: MAN

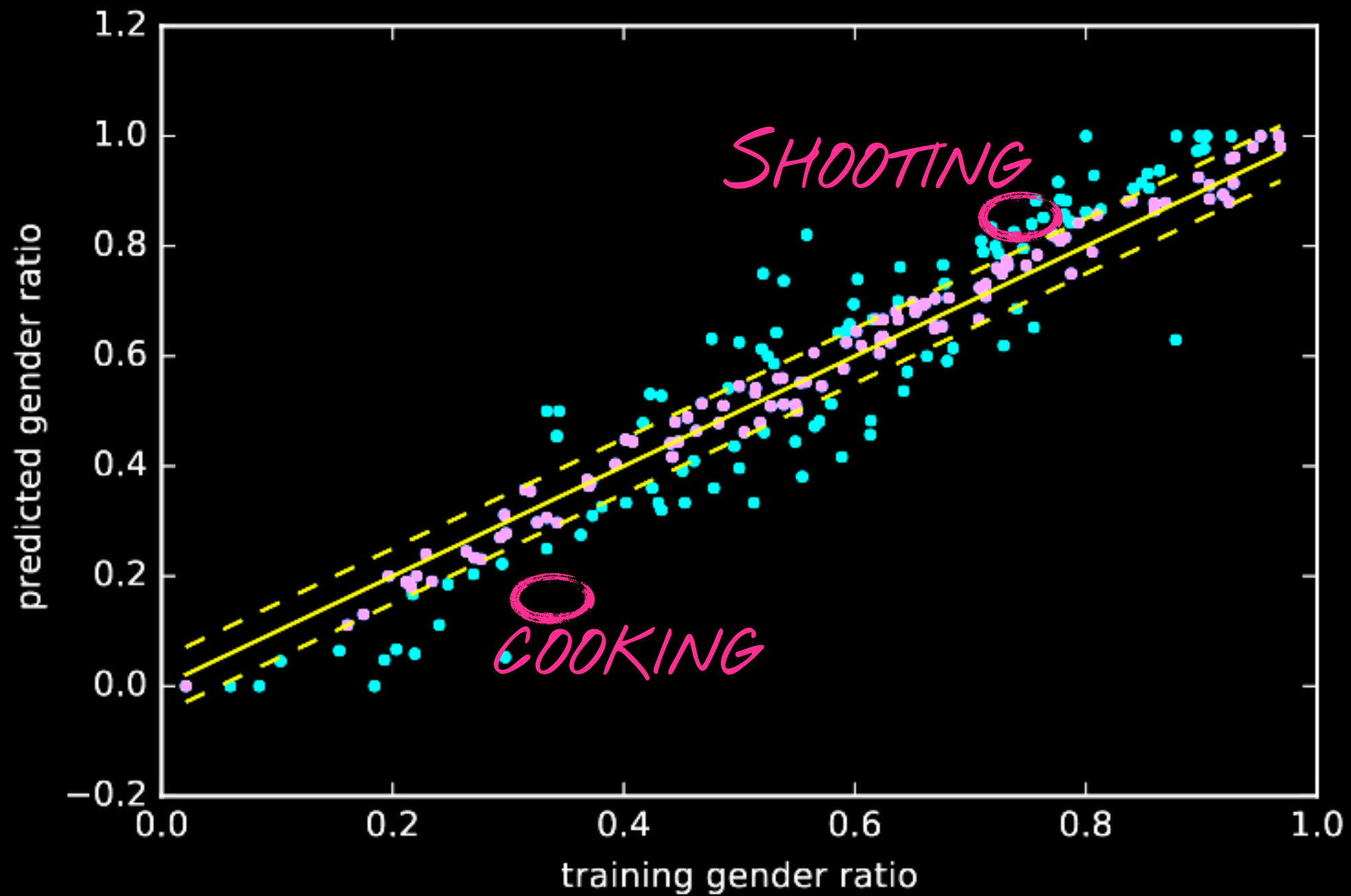


Agent: WOMAN

ILP



Results





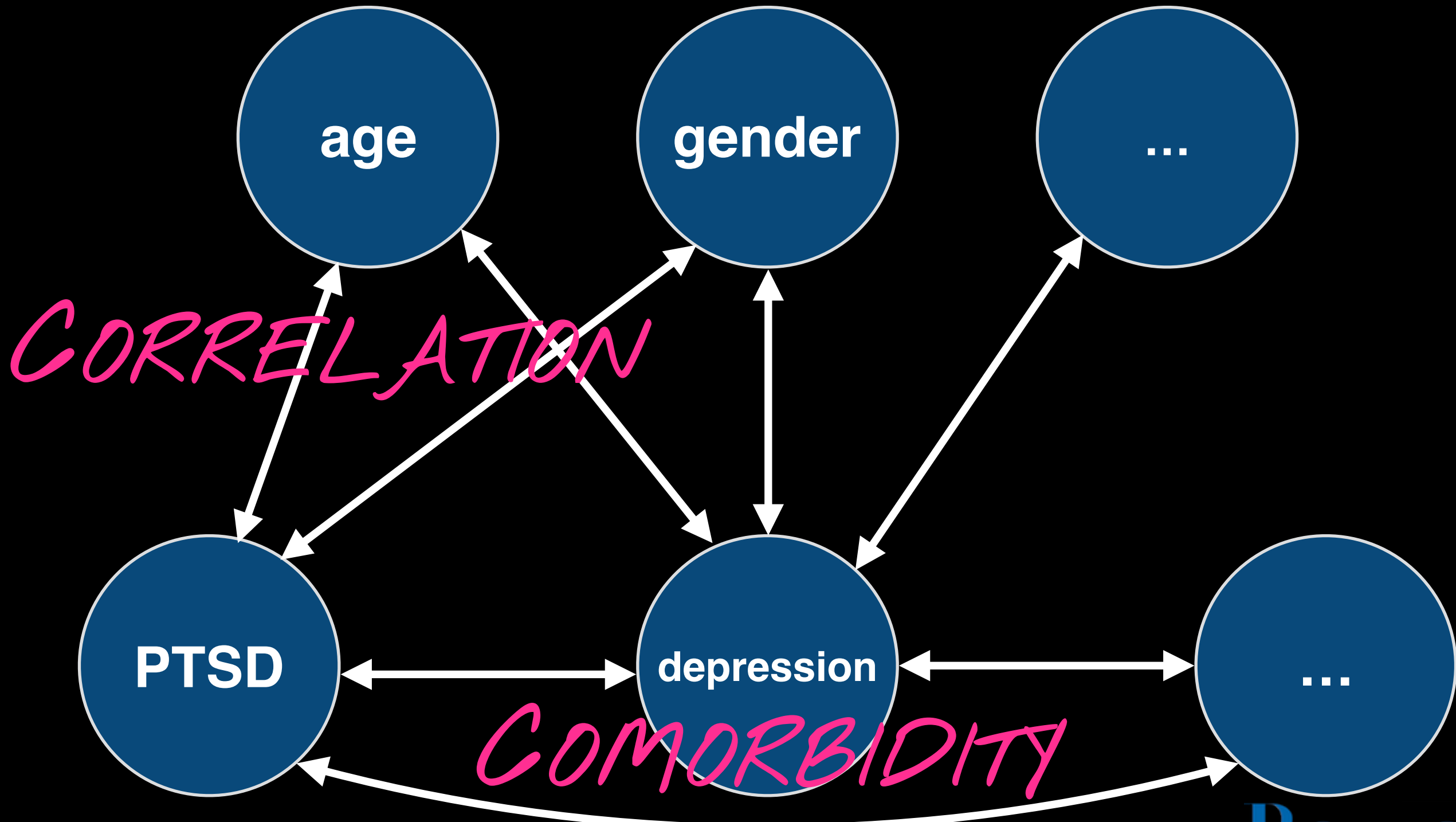
Idea!

*ADD DEMOGRAPHIC
COMPONENT IN MODEL*

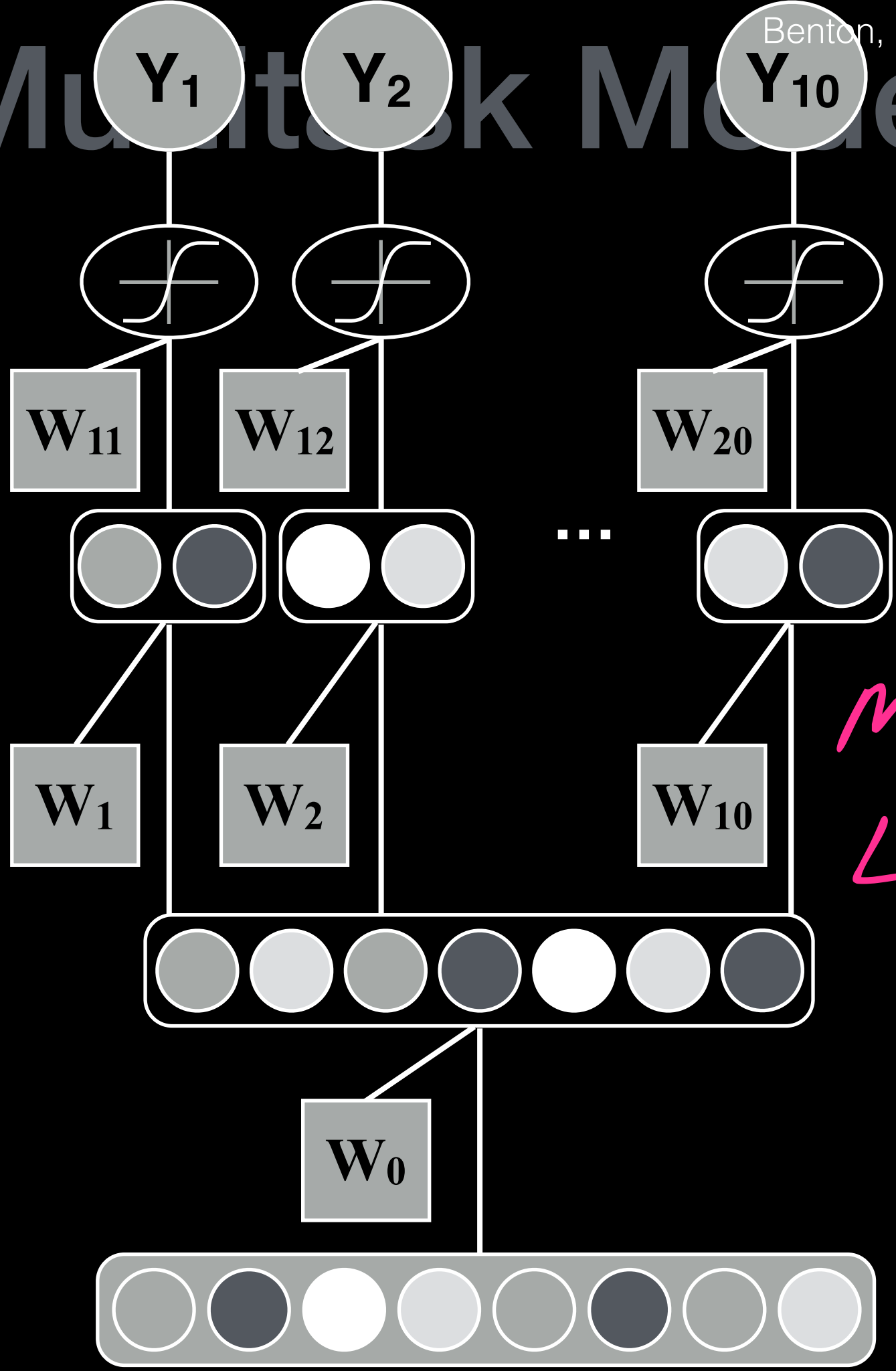




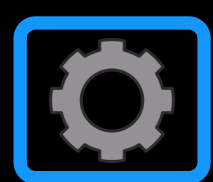
Comorbidity and Correlation



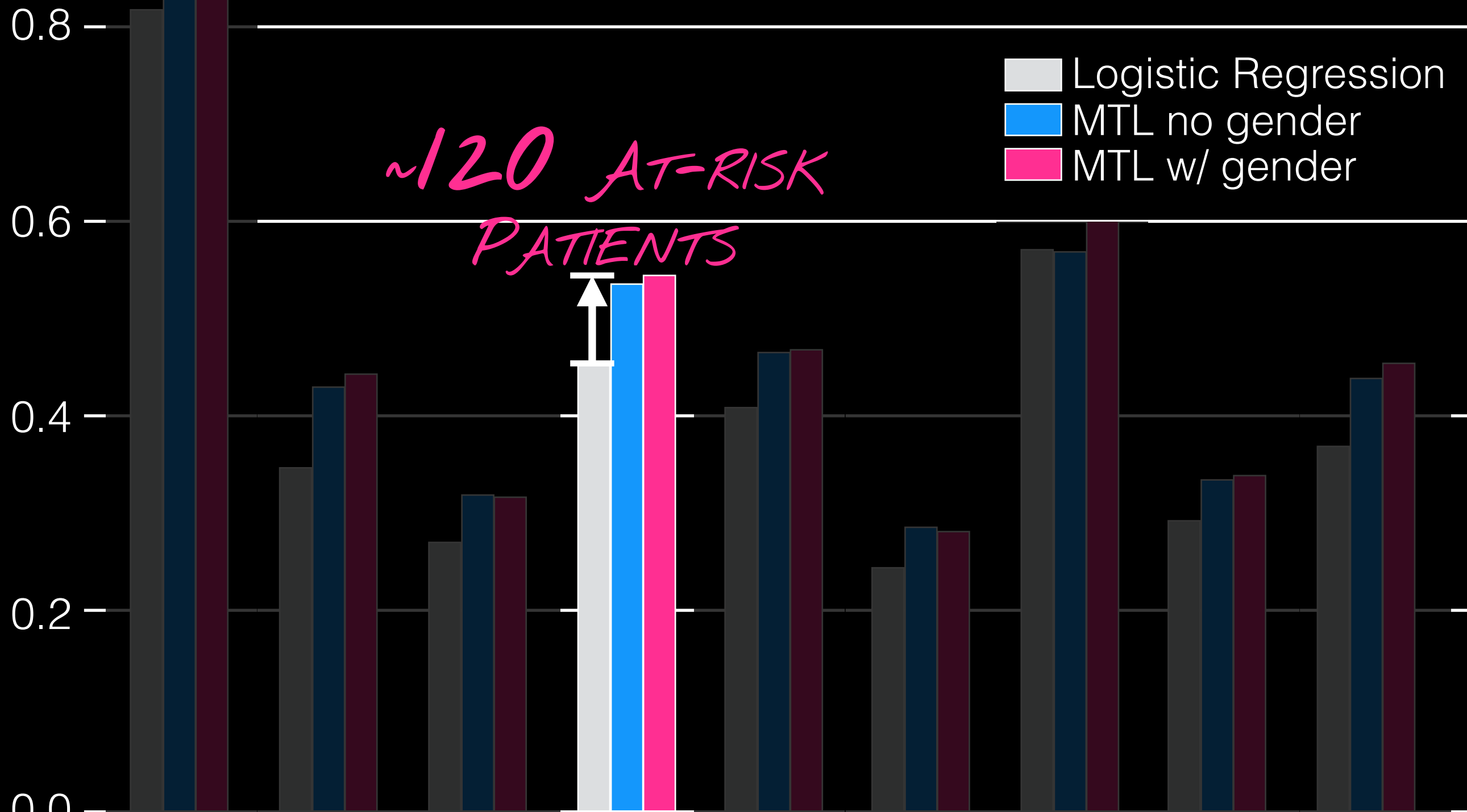
Multitask Model



*MULTITASK
LEARNING*



Results: TPR@FPR=0.1



NNT | Depression | Anxiety | Suicide | Eating | Schizophrenia | Panic | PTSD | Bipolar

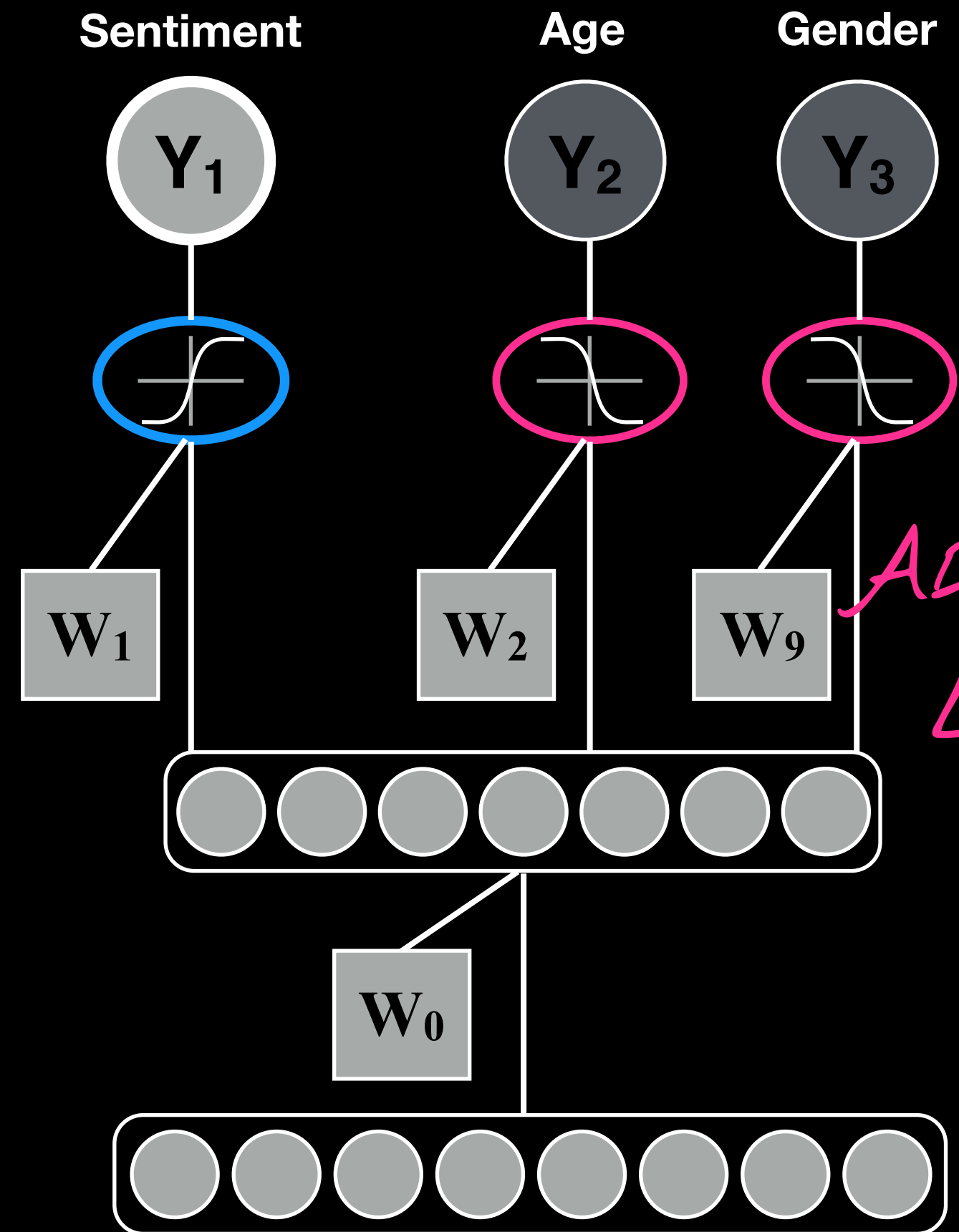
4791 | 2407 | 1400 | 1208 | 749 | 349 | 263 | 248 | 191

Idea!

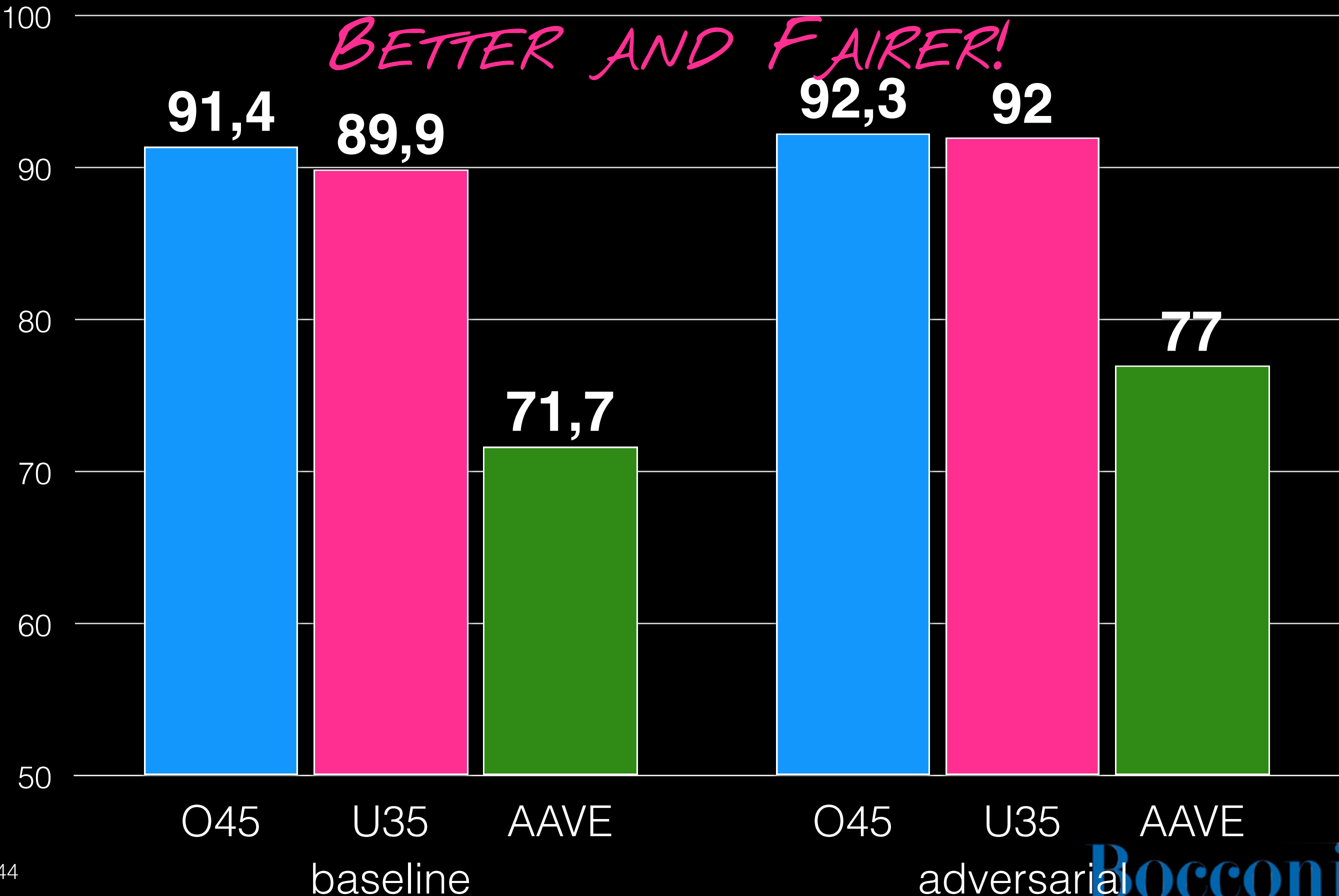
CORRECT FOR BIAS ADVERSARIALLY



Adversarial Model

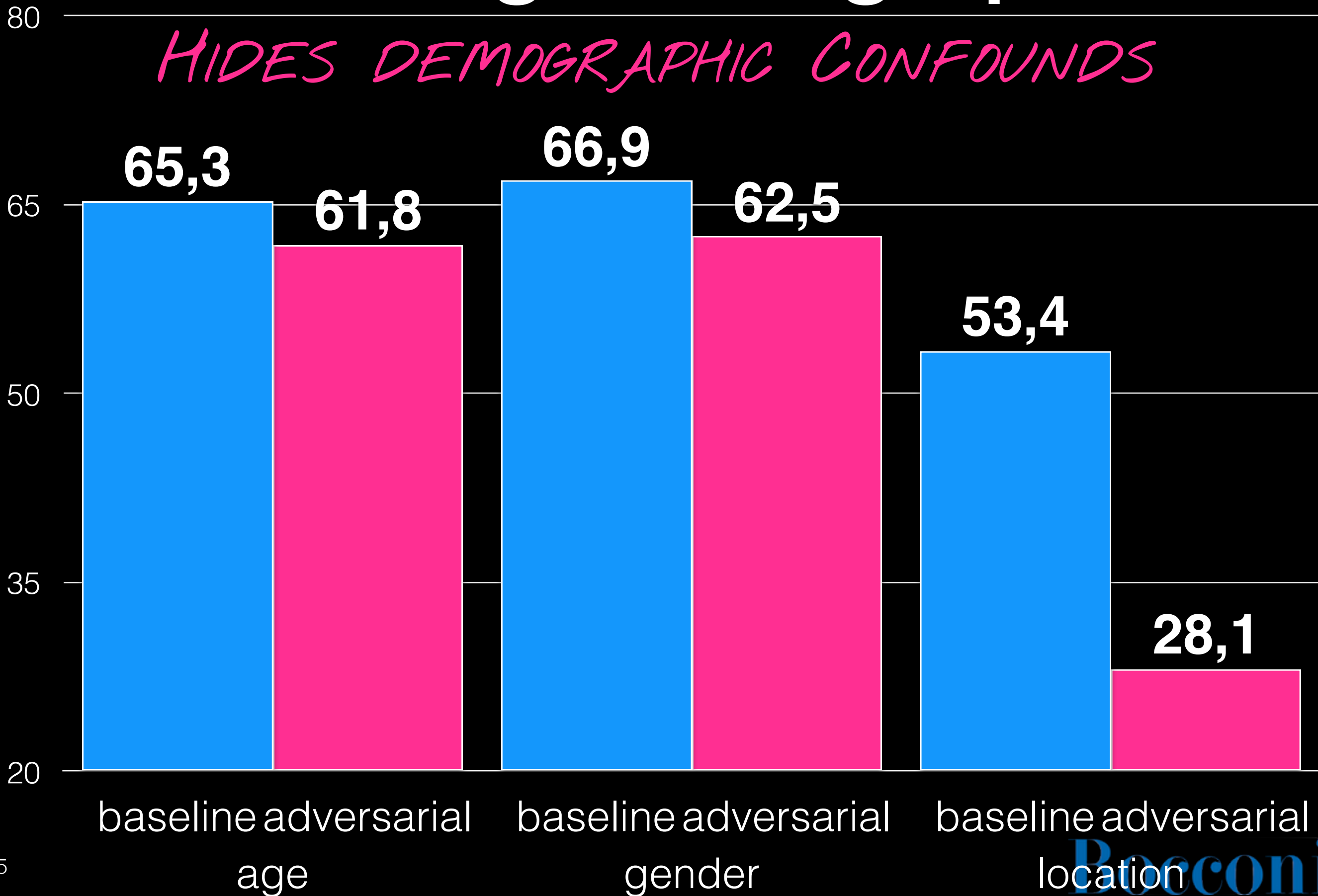


Results



Protecting Demographics

HIDES DEMOGRAPHIC CONFOUNDS

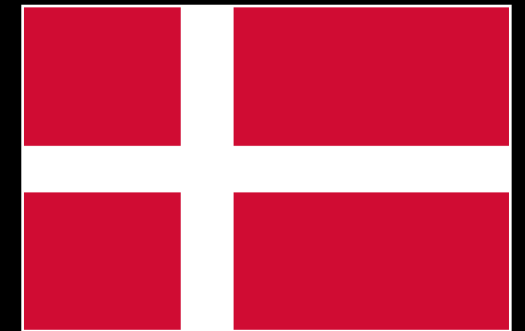
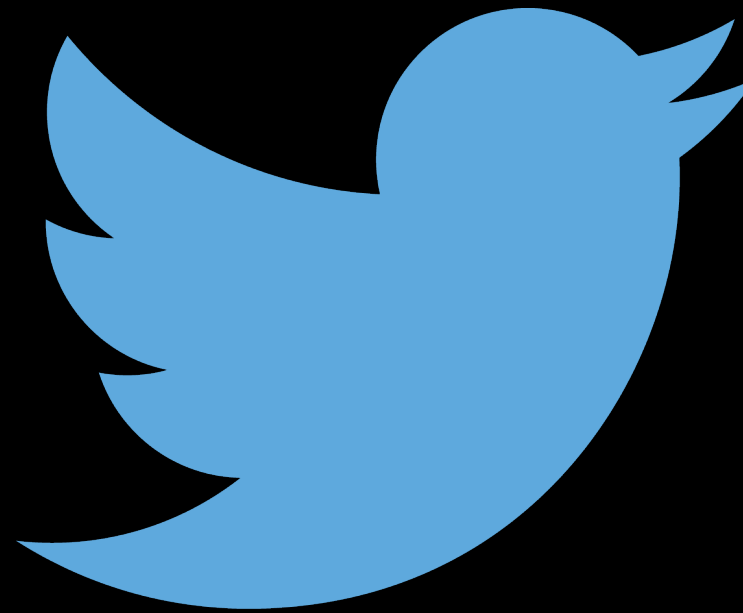
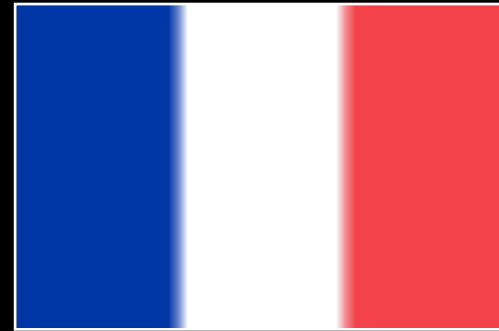
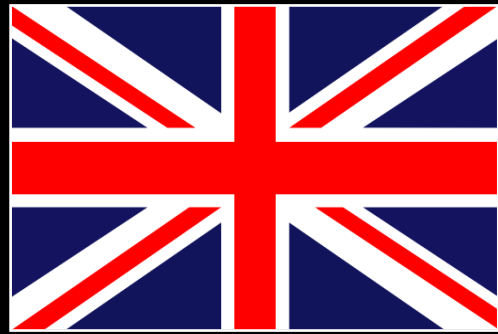


Part 4: Design Bias



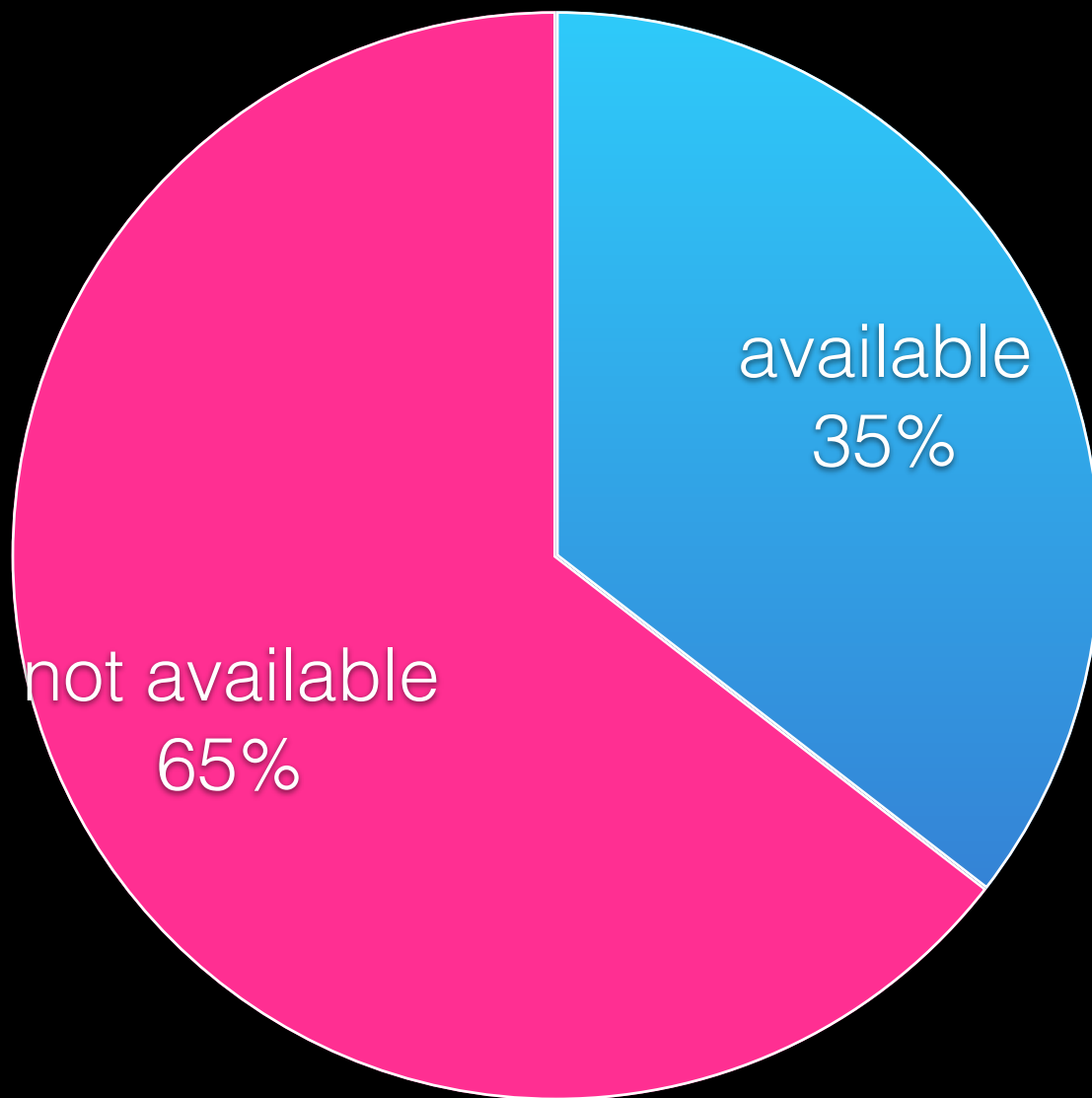


Exposure

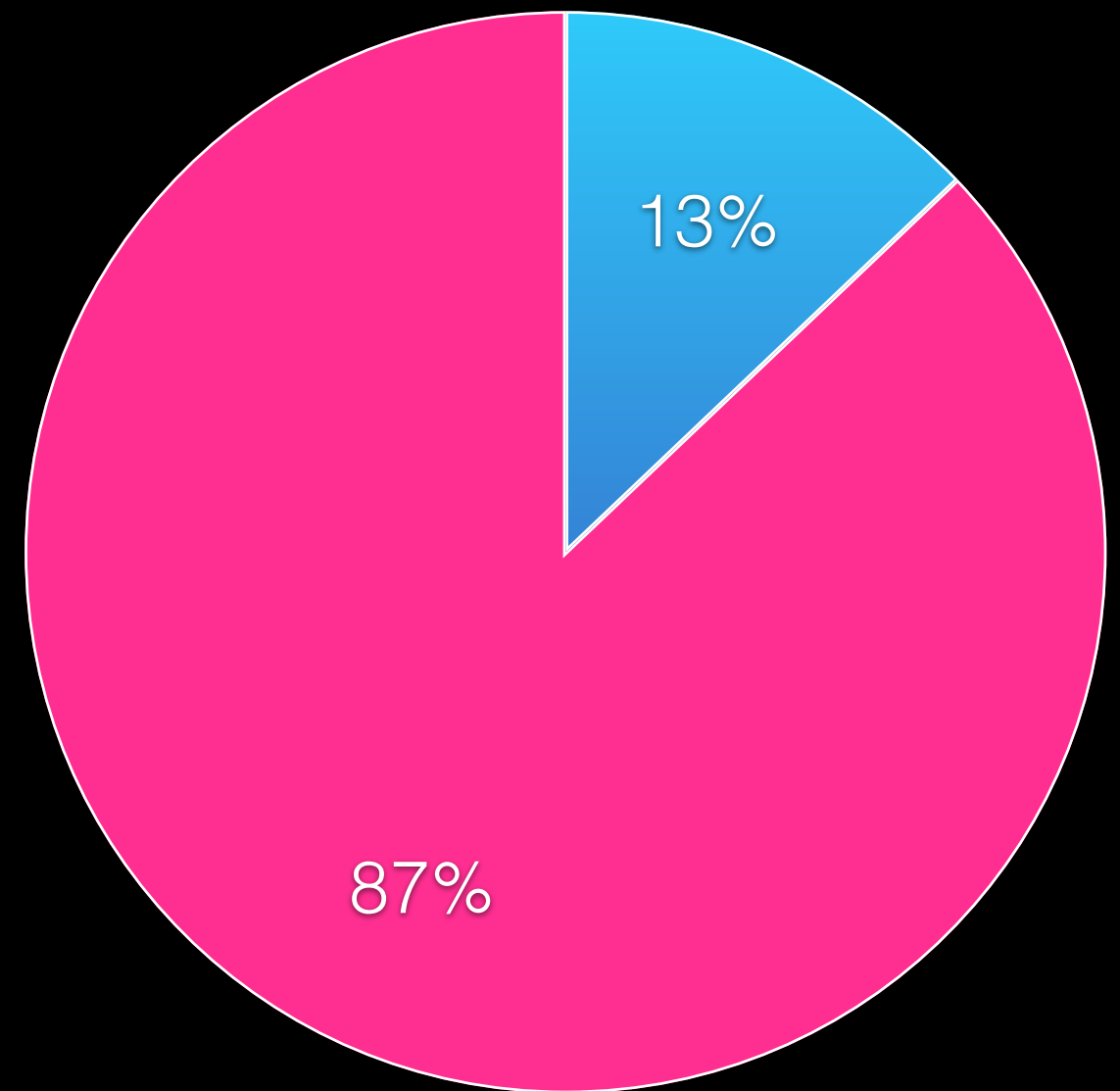


Under-Exposure

treebanks *



semantic resources



**BEFORE UD...* evaluation

Over-Exposure



Aerial view of the New York City skyline with the Freedom Tower. The text 'American English' is overlaid in large white font, with 'New York City' and '85m' in smaller white font below it.

American
New York City
English
85m



Aerial view of the Lagos skyline across a body of water with several sailboats. The text 'Nigerian English' is overlaid in large white font, with '16m' in smaller white font below it.

Nigerian
English
16m

POS tagging

Discourse



Dual Use

Task	Pro	Con
authorship attribution	historical documents	dissenter anonymity
text classification	sentiment analysis	censorship
personalization	better user experience	tailored ads

Normative vs Descriptive Ethics

The image shows two screenshots of the Google Translate interface. The top screenshot shows the English input "She is a doctor. He is a nurse." being translated into Turkish as "O bir doktor. O bir hemşire." The bottom screenshot shows the Turkish input "O bir doktor. O bir hemşire" being translated back into English as "He is a doctor. She is a nurse".

NORMATIVELY WRONG

DESCRIPTIVELY WRONG

Normative vs Descriptive Ethics



why are american

why are american **so fat**

why are american **so long**

why are american **so proud**

why are american **houses made of wood**

why are american **trucks different to european**

why are american **universities the best**

why are american **houses made of cardboard**

why are american **cars so big**

Google-Suche

Auf gut Glück!

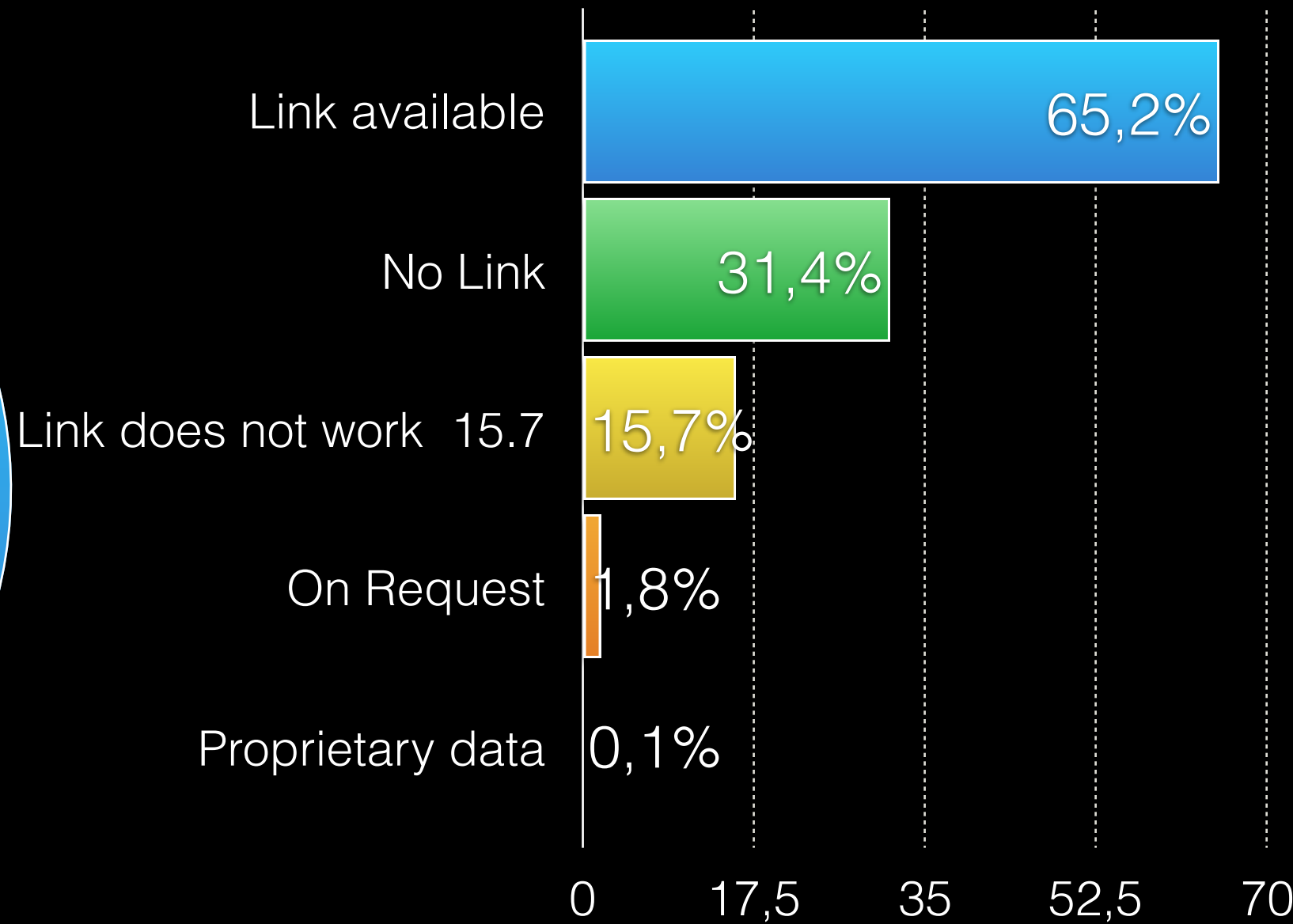
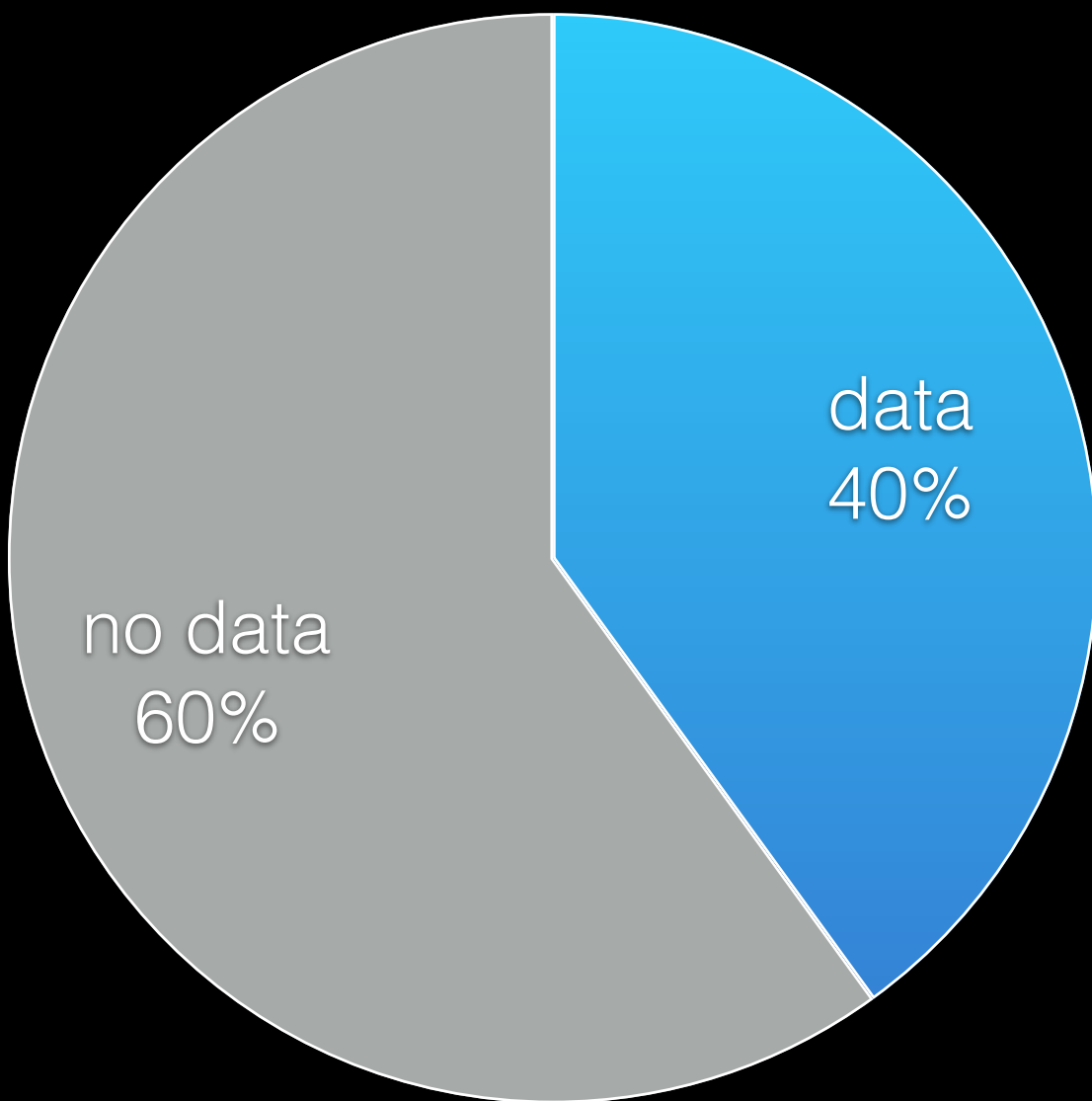
[Weitere Informationen](#)

[Unangemessene Vervollständigungen melden](#)

NORMATIVELY WRONG

DESCRIPTIVELY TRUE?

Replicability: Data



Replicability: Significance

Cut-offs: 0.1 (meh), 0.05 (standard), 0.01 (strict)

<p>(barely) not statistically significant (p=0.052) a barely detectable statistically significant difference (p=0.073) a borderline significant trend (p=0.09) a certain trend toward significance (p=0.08) a clear tendency to significance (p=0.052) a clear trend (p<0.09) a clear, strong trend (p=0.09) a considerable trend toward significance (p=0.069) a decreasing trend (p=0.09) a definite trend (p=0.08) a distinct trend toward significance (p=0.07) \borderline conventional significance (p=0.051) borderline level of statistical significance (p=0.053)</p>	<p>borderline significant (p=0.09) did not quite reach conventional levels of statistical significance (p=0.079) did not quite reach statistical significance (p=0.063) did not reach the traditional level of significance (p=0.10) did not reach the usually accepted level of clinical significance (p=0.07) difference was apparent (p=0.07) direction heading towards significance (p=0.10) does not appear to be sufficiently significant (p>0.05) does not narrowly reach statistical significance (p=0.06)</p>	<p>does not reach the conventional significance level (p=0.098) effectively significant (p=0.051) equivocal significance (p=0.06) essentially significant (p=0.10) extremely close to significance (p=0.07) failed to reach significance on this occasion (p=0.09) failed to reach statistical significance (p=0.06) fairly close to significance (p=0.065) fairly significant (p=0.09) falls just short of standard levels of statistical significance (p=0.06) fell (just) short of significance (p=0.08)</p>	<p>fell barely short of significance (p=0.08) scarcely significant (0.05<p>0.1) significant at the .07 level significant tendency (p=0.09) significant to some degree (0<p>1) significant, or close to significant effects (p=0.08, p=0.05) significantly better overall (p=0.051) significantly significant (p=0.065) similar but not nonsignificant trends (p>0.05) slight evidence of significance (0.1>p>0.05) slight non-significance (p=0.06) slight significance (p=0.128)</p>
---	---	---	---



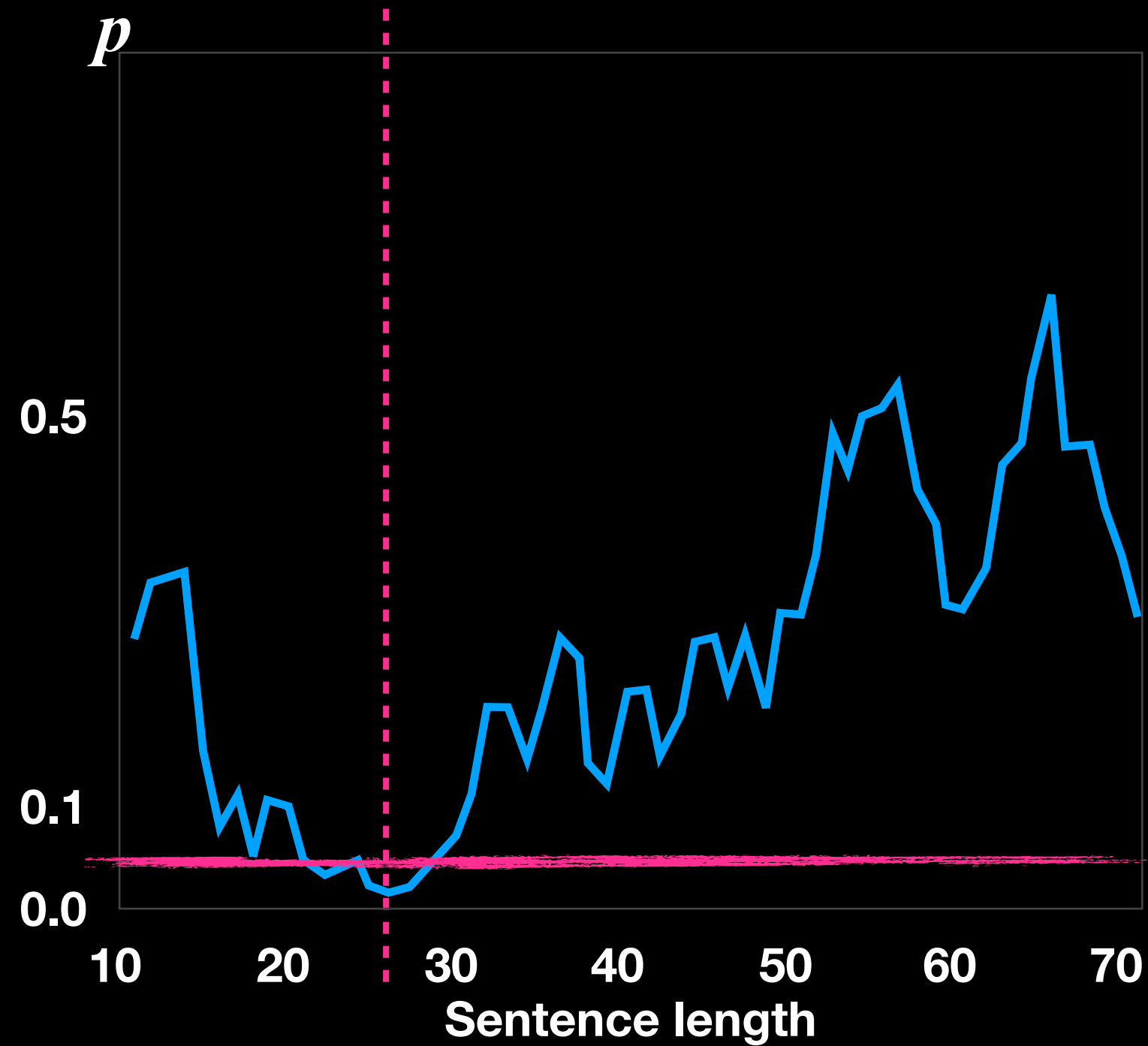
Don't choose among metrics

metric	p
precision	0,0899
precision	0,062
recall	0,179
accuracy	0,0014

REPORT!



Don't choose sample sizes



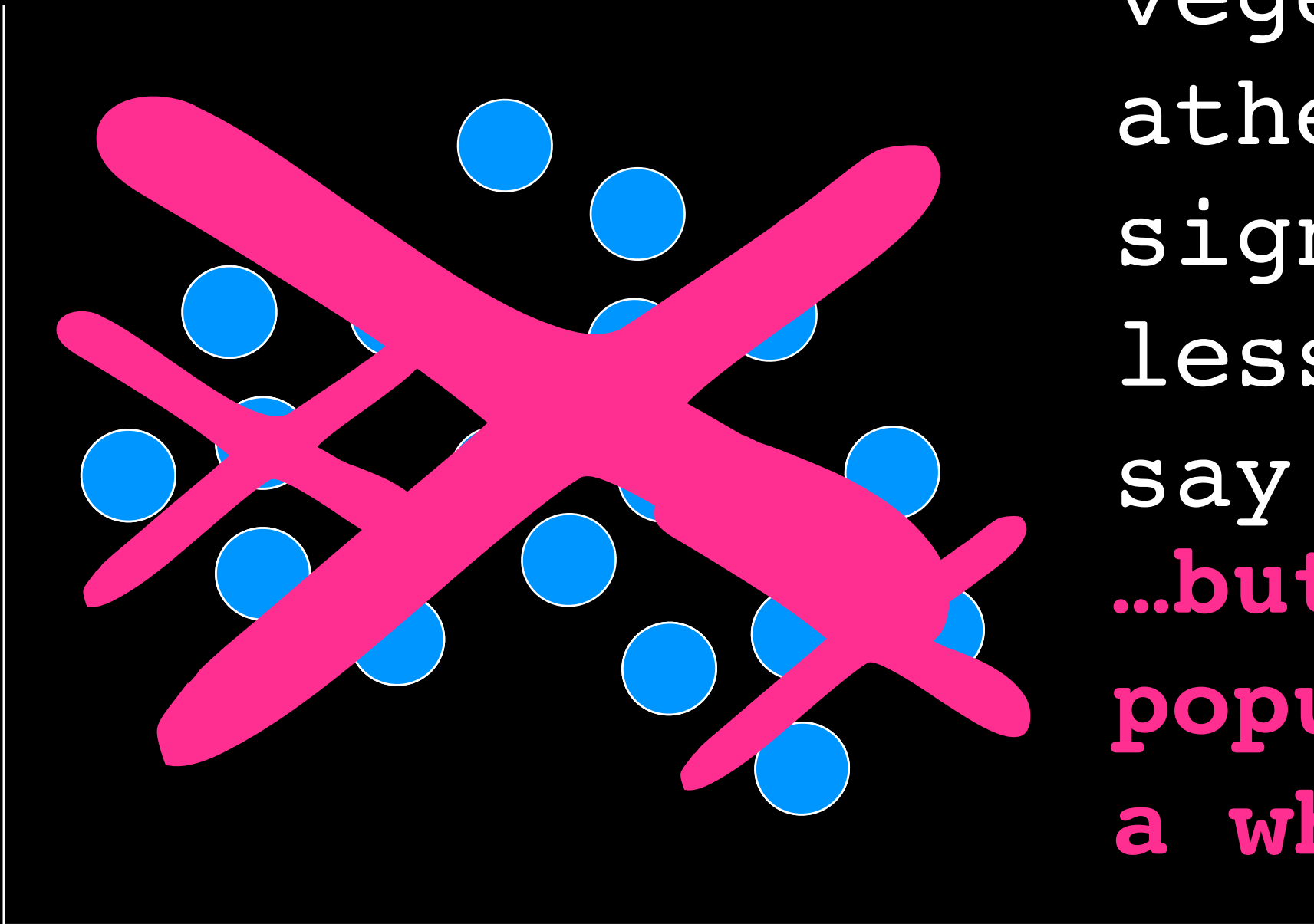
"We observed significant results at a sentence length of 26"
...but not with smaller or larger sentences!



Don't Choose Subsets

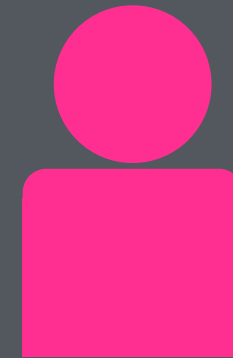
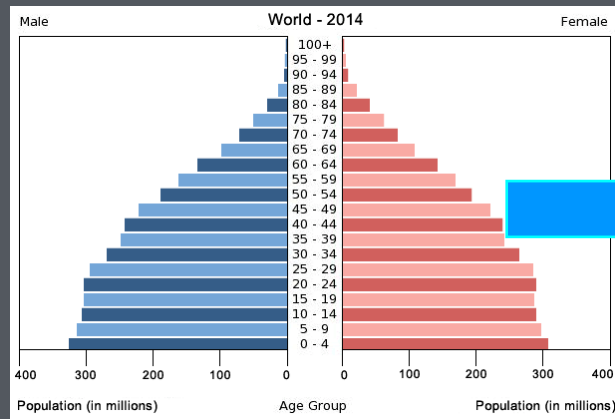
"Young, left-handed, vegetarian atheists are significantly less likely to say X"

...but the population as a whole isn't!



Wrapping Up

Sources of Bias



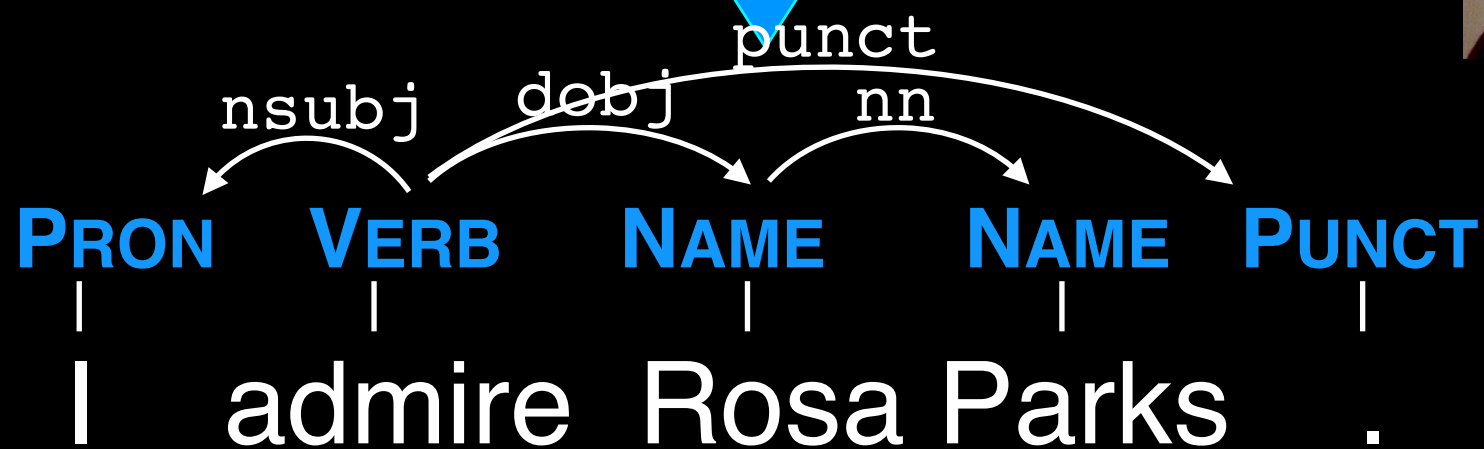
SELECTION

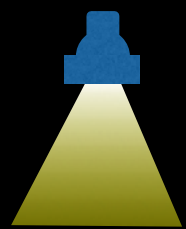
ANNOTATION

MODELS



DESIGN



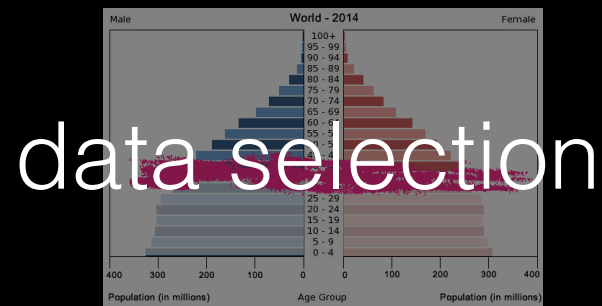


What can we do?

Source

Problem

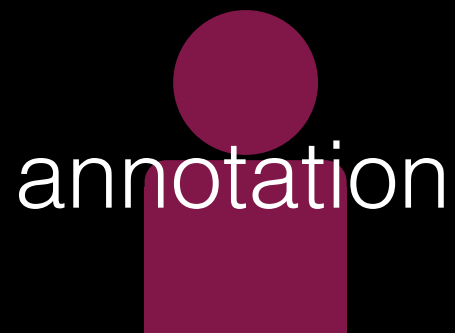
Countermeasures



data selection

Exclusion

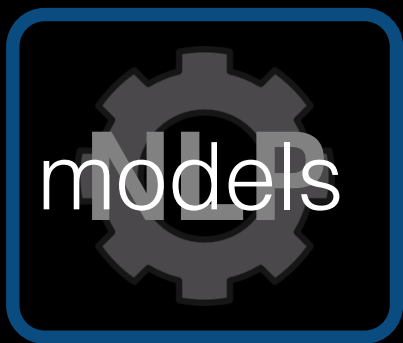
stratification, priors



annotation

Label Bias

annotation models,
disagreement weighting



models

Overgeneralization

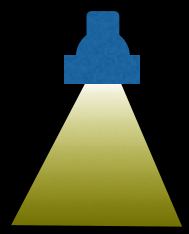
dummy labels, error weighting,
adversarial learning



research
design

Exposure

always consider possible
impact



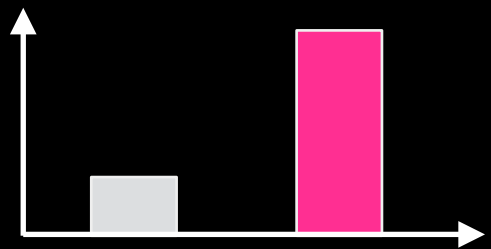
The Goals



Fairness



Personalization



Performance

Take-home points

- Beware of **bias** from **data**, **models**, and **design**
- Apply **countermeasures** and check
- Ask yourself:
"Am I comfortable with my system classifying me?"

www.dirkhovy.com/portfolio/papers

Thank you!



@dirk_hovy

www.dirkhovy.com

Bocconi

www.dirkhovy.com/portfolio/papers

Questions?



@dirk_hovy

www.dirkhovy.com

Bocconi