



Five Sources of Biases and Ethical Issues in NLP, and What to Do about Them

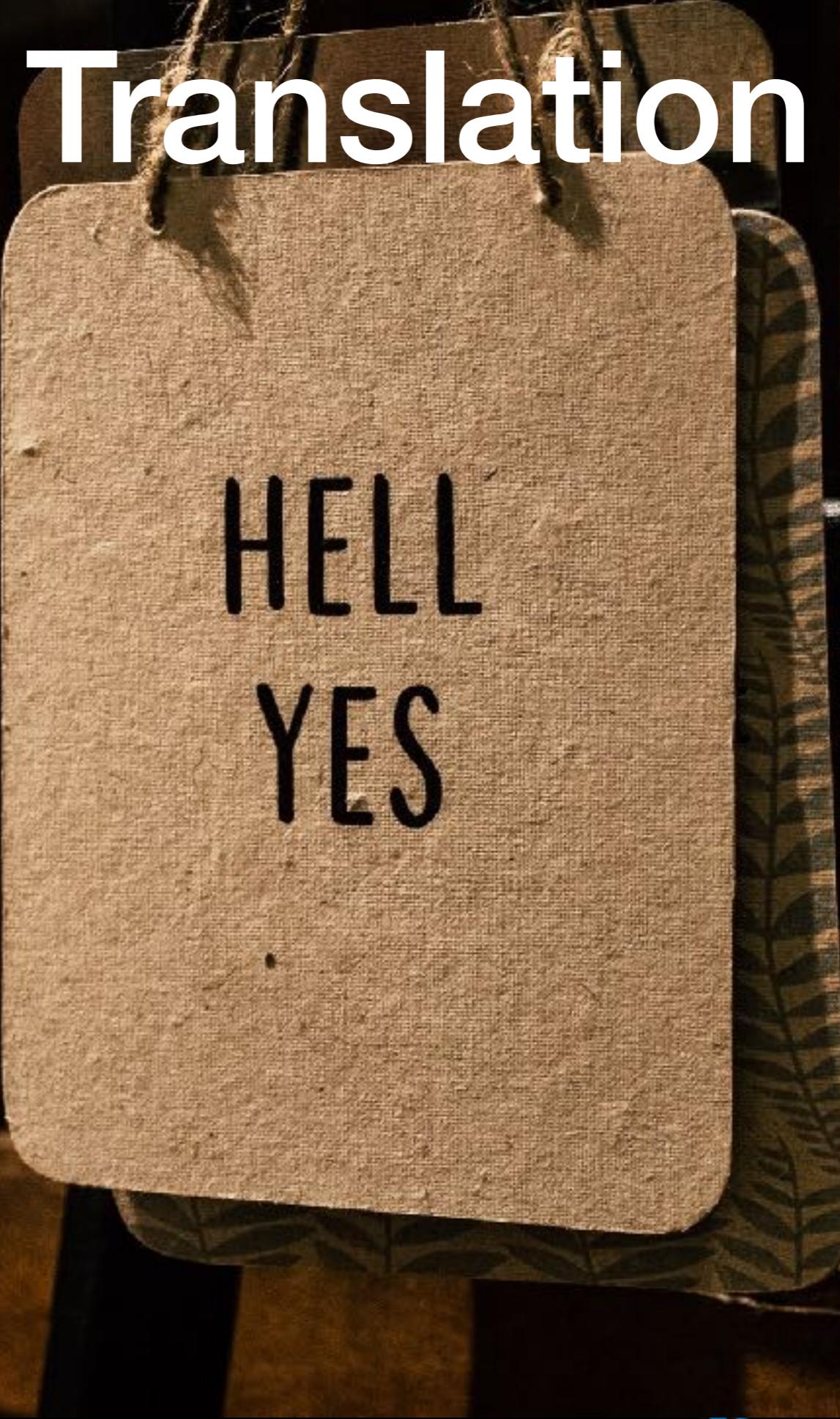
Dirk Hovy
Bocconi University, Milan

www.dirkhovy.com
dirk.hovy@unibocconi.it

 @dirk_hovy

Bocconi

Machine Translation



Bocconi

Text Generation



In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

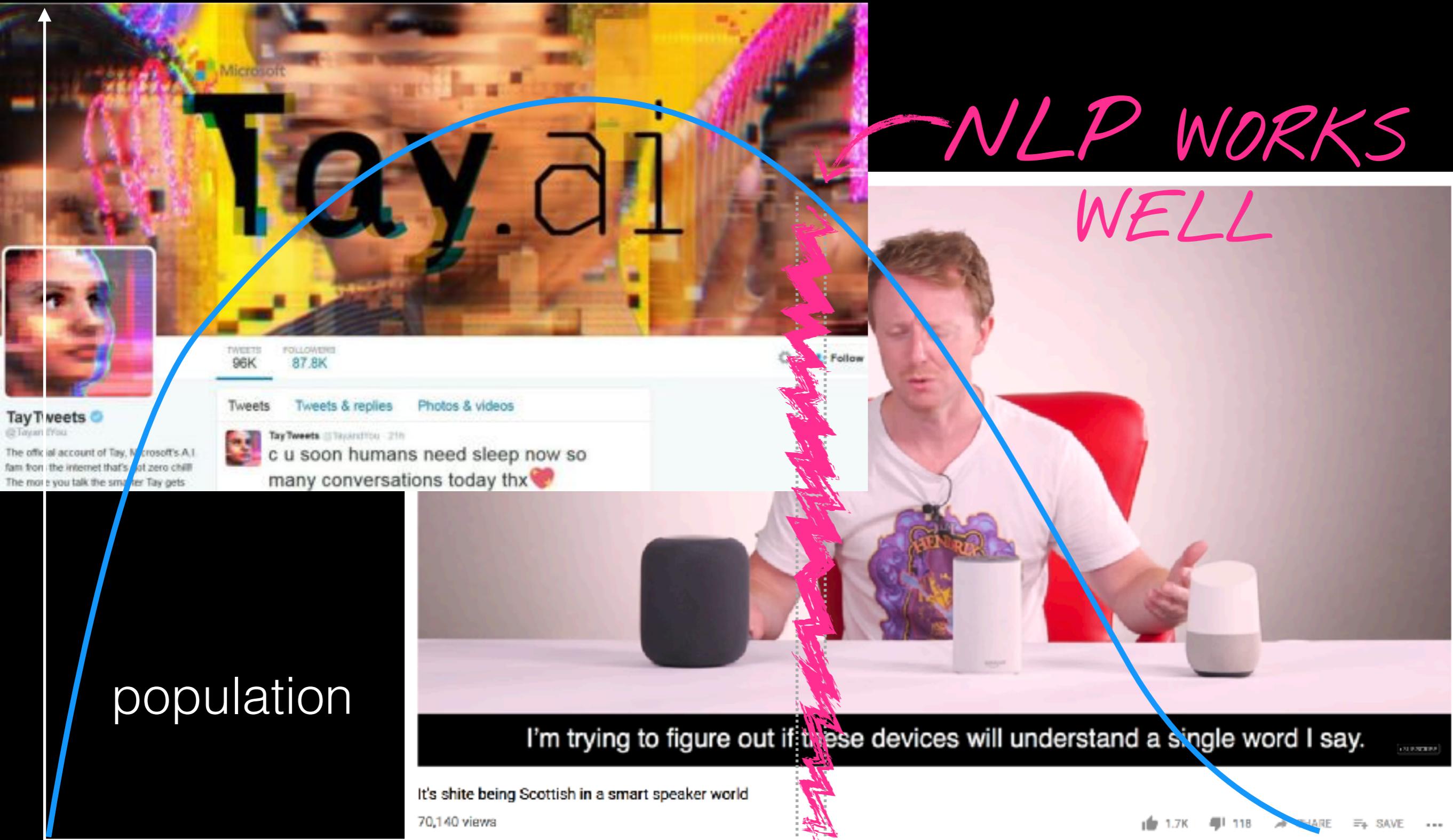
The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Biased Systems



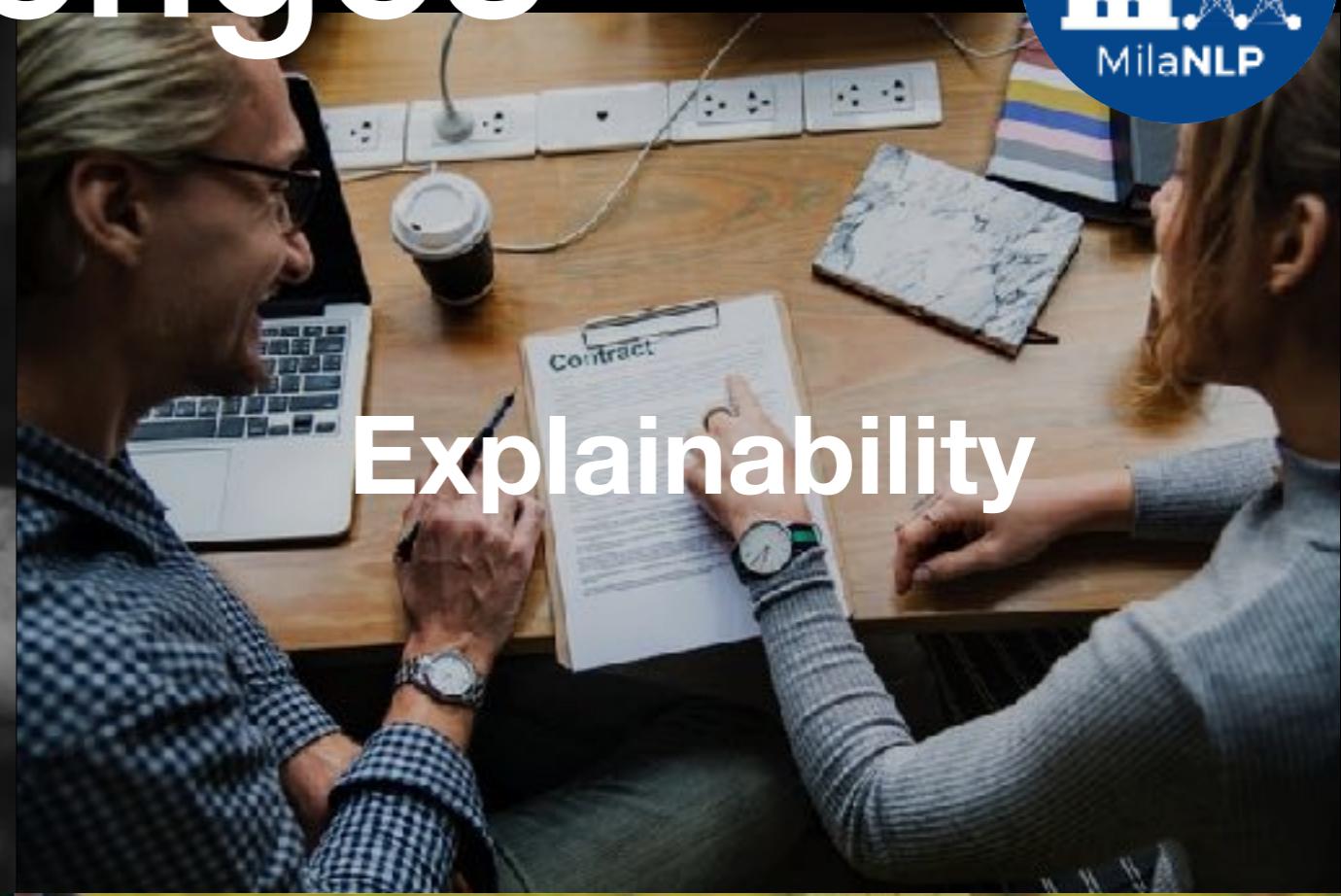
Challenges



Trust



Explainability



Recourse



Fairness



Bocconi



Goals for Today

- Introduce **terminology** and **concepts** of **bias**
- Show five **sources** of bias in NLP
- Discuss **consequences**
- Introduce **countermeasures**
- For a paper, see Shah, Schwartz & Hovy (ACL 2020):
<https://www.aclweb.org/anthology/2020.acl-main.468v2.pdf>



Terminology

Bocconi

Bias

Bias ≠ Bad

Ethics ≠ Bias

BIAS!

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Normative vs Descriptive Ethics



English Turkish Spanish Detect language ▾

English Turkish Spanish ▾

Translate

She is a doctor.
He is a nurse.

O bir doktor.
O bir hemşire.

🔊 🔍 ⌨ ▾

31/5000

☆ 🔍 🔊 <

English Turkish Spanish

Turkish - detected

English Turkish Spanish ▾

Translate

O bir doktor.
O bir hemşire

🔊 🔍

28/5000

He is a doctor.
She is a nurse ✅

☆ 🔍 🔊 <

NORMATIVELY WRONG
DESCRIPTIVELY WRONG

Bocconi

Normative vs Descriptive Ethics



why are american

why are american **so fat**

why are american ~~so~~ 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

NORMATIVELY WRONG
DESCRIPTIVELY TRUE?

Dual Use



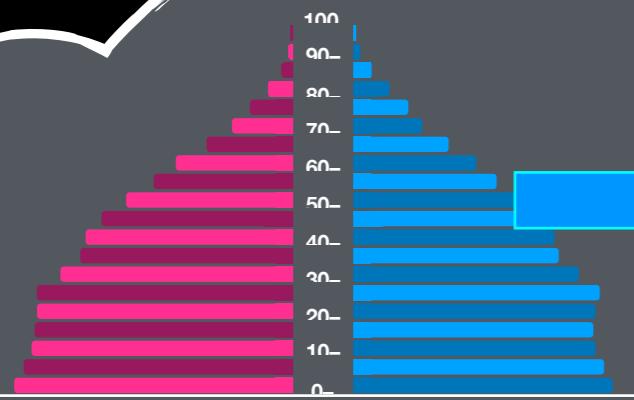
+ INTENDED USE

- UNINTENDED USE
OR CONSEQUENCES



Sources of Bias

Sources of Bias



DATA



ANNOTATION



REPRESENTATIONS

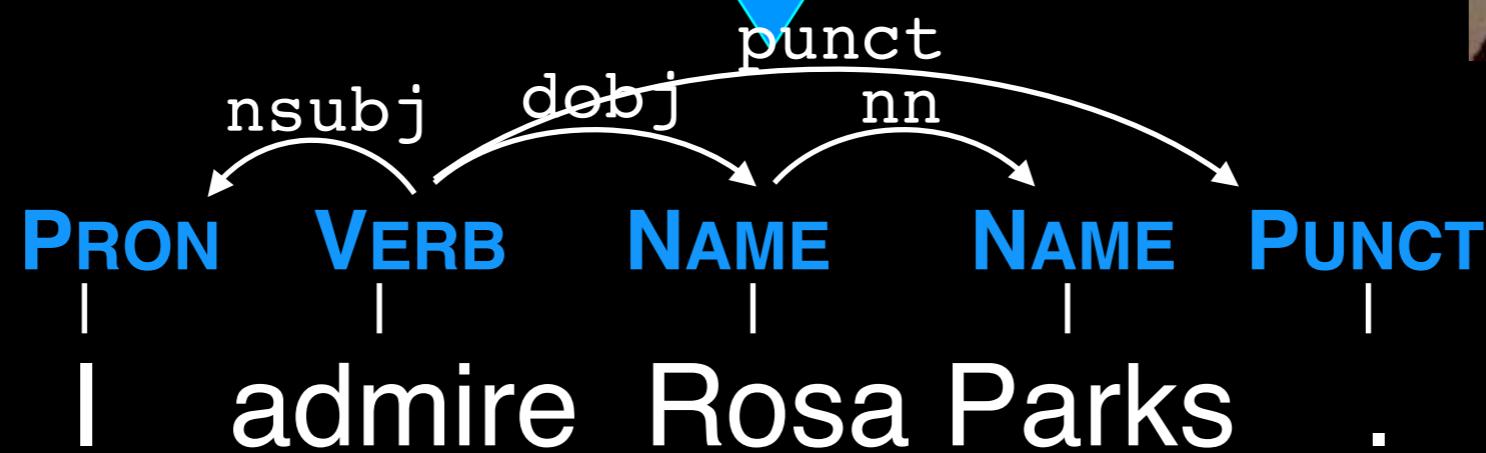


MODELS

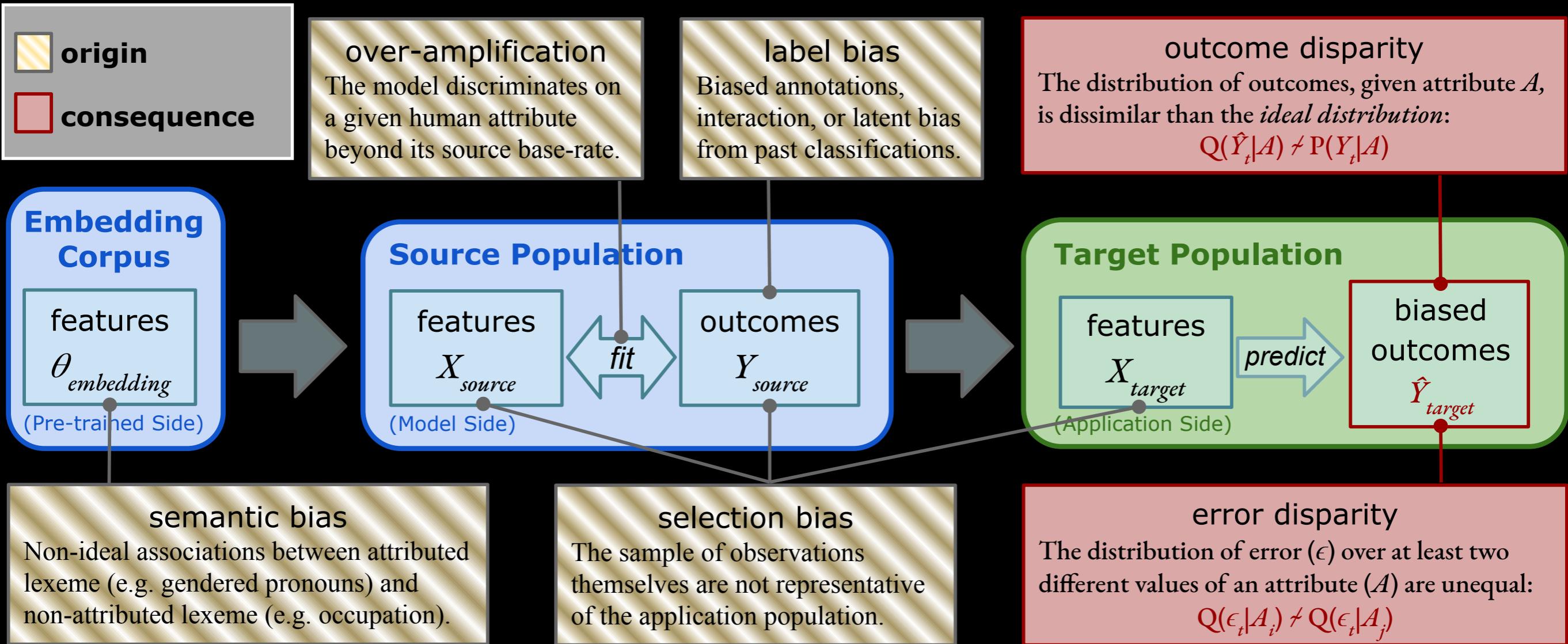
NLP



DESIGN

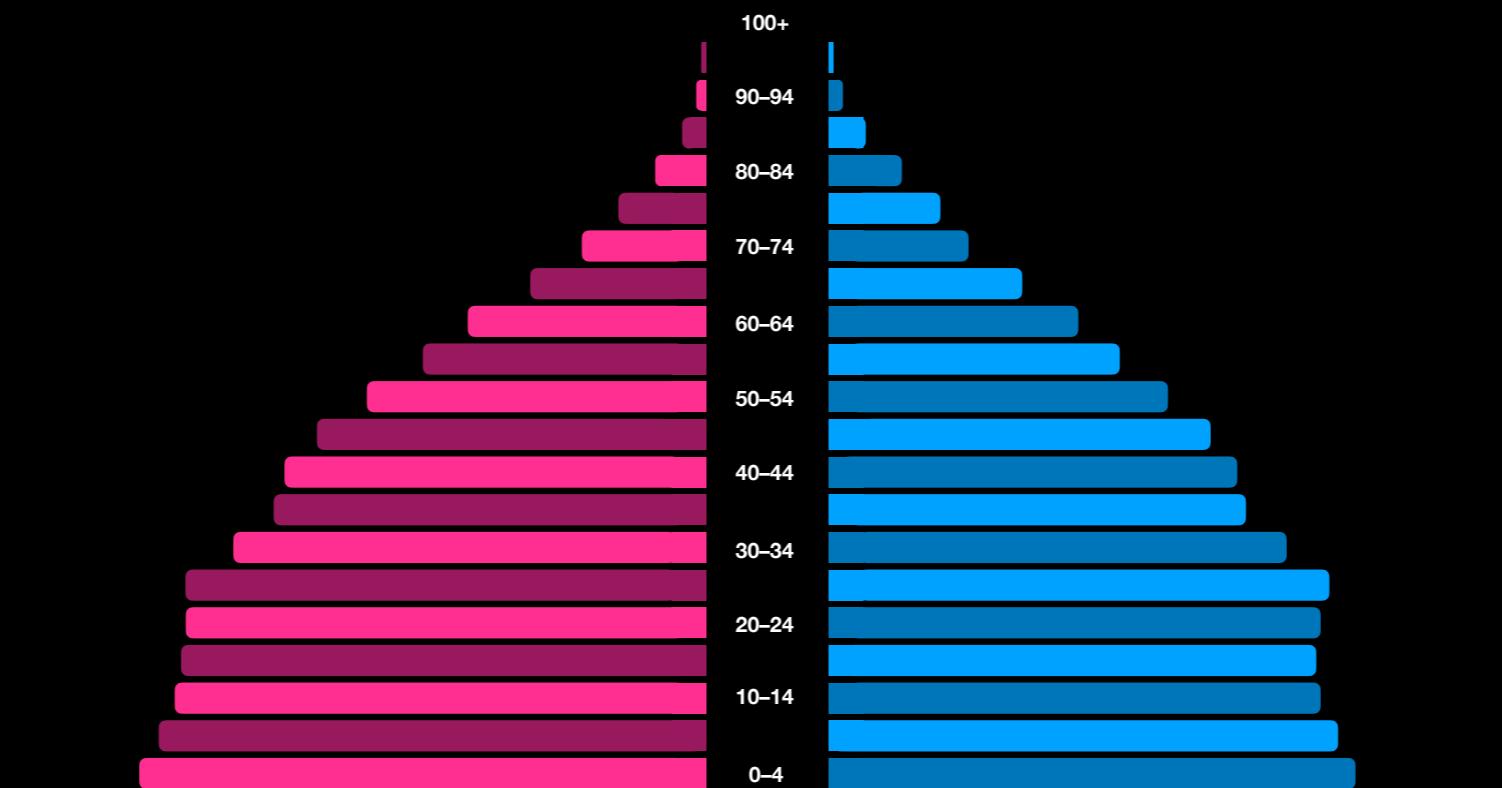


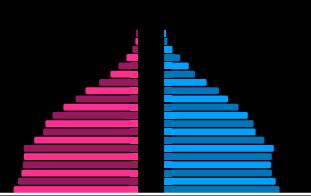
Formally...



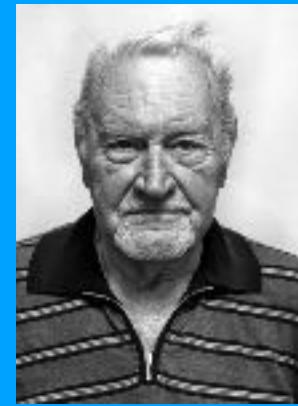
Part 1:

Selection Bias





Language Varies

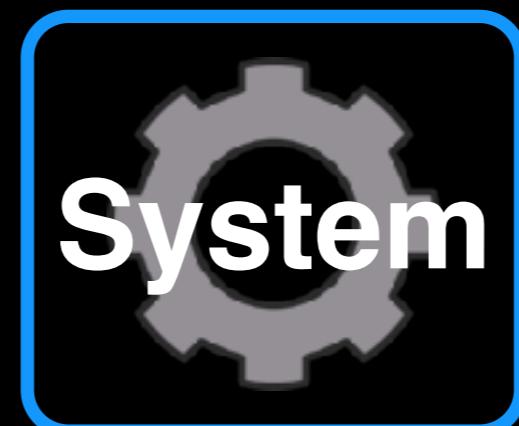


Example 1

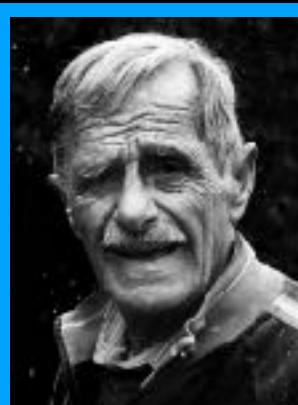
I don't understand you...



Example 2



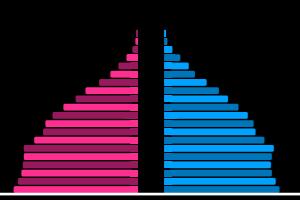
Hello,
computer



Example N

Shite...



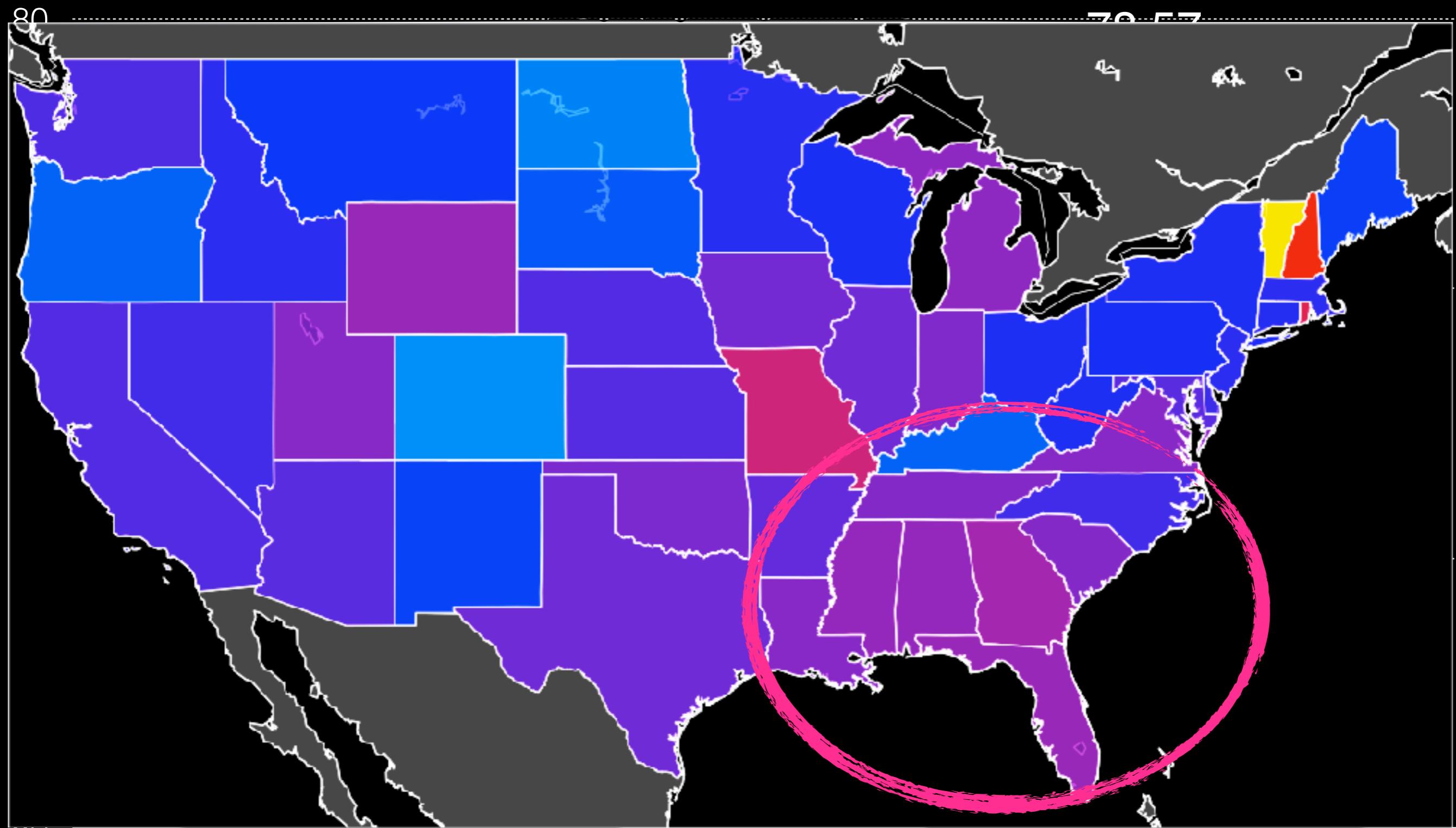


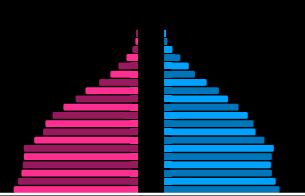
Jørgensen et al. (WNUT 2015)
Hovy & Spruit (ACL 2016)



Exclusion

F1





Exclusion

Hovy & Søgaard (ACL 2015)
Hovy & Spruit (ACL 2016)



accuracy

100

95

90

85

80

100+

90–94

80–84

70–74

60–64

50–54

40–44

30–34

20–24

10–14

0–4

POS-tagging

600 user reviews

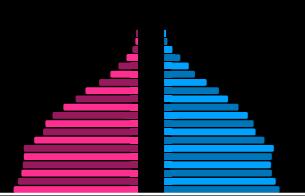
CoNLL trained

O45

U35

O45

U35
Bocconi



Wrong Coreference

Mention -----coref----- Mention -----coref----- Mention -----coref----- Mention
The surgeon could n't operate on his patient : it was his son !

Mention -----coref----- Mention -----coref----- Mention -----coref----- Mention
The surgeon could n't operate on their patient : it was their son !

Mention -----coref----- Mention -----coref----- Mention -----coref----- Mention
The surgeon could n't operate on her patient : it was her son !

More generally...

The sample of observations is not representative of the application population

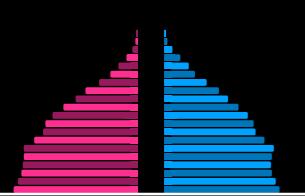
DISTRIBUTION
OF ATTRIBUTE
IN TRAINING DATA

$$Q(A_S)$$

\neq

$$P(A_t)$$

IDEAL DISTRO
OF ATTRIBUTE
IN TARGET DATA



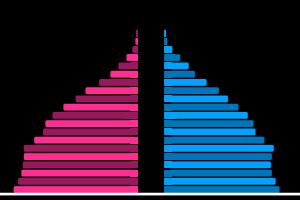
Idea!

Volkova et al. (2013)
Hovy (2015)
Lynn et al. (2017)



INCLUDE DEMOGRAPHIC INFORMATION
IN TEXT REPRESENTATION



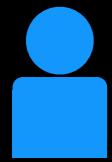


Systems

AGNOSTIC



training



This is a tiny little example text written by someone.

data



training



This is a tiny little example text written by someone.

data



training

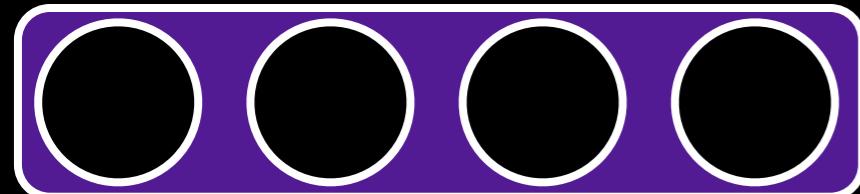


This is a tiny little example text written by someone.

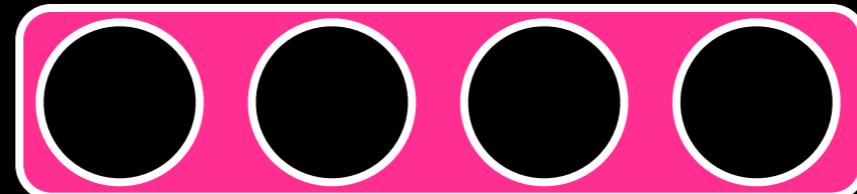
data

This is a tiny little example text written by someone.

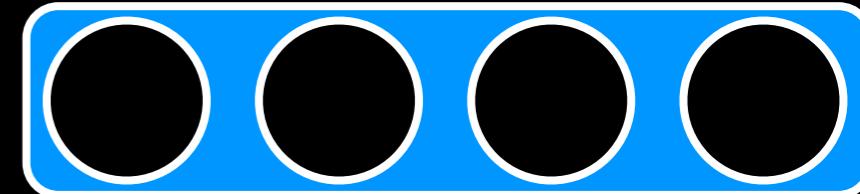
+

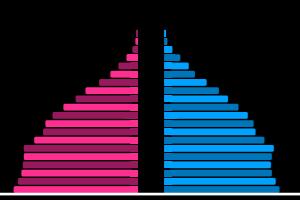


+



+





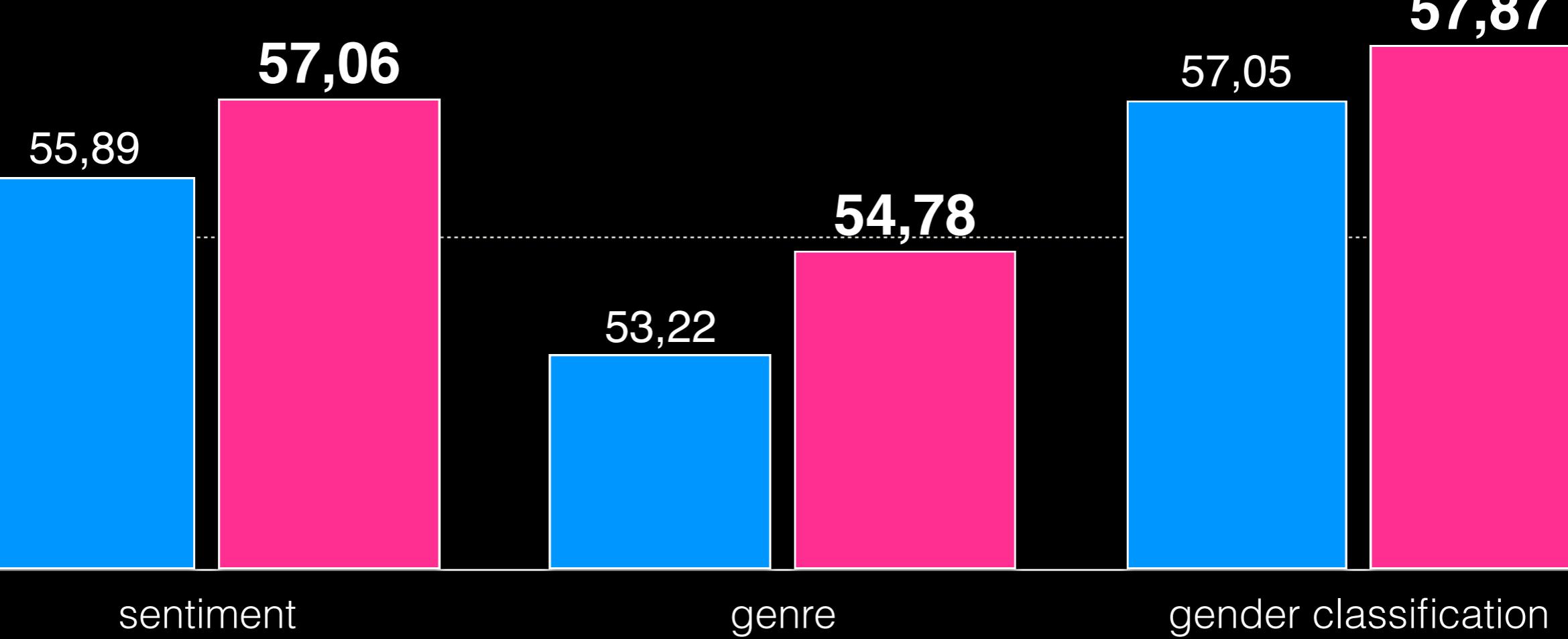
Results for Age (avg)

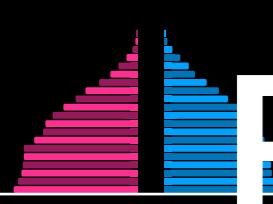
F1

65

- agnostic
- aware

60



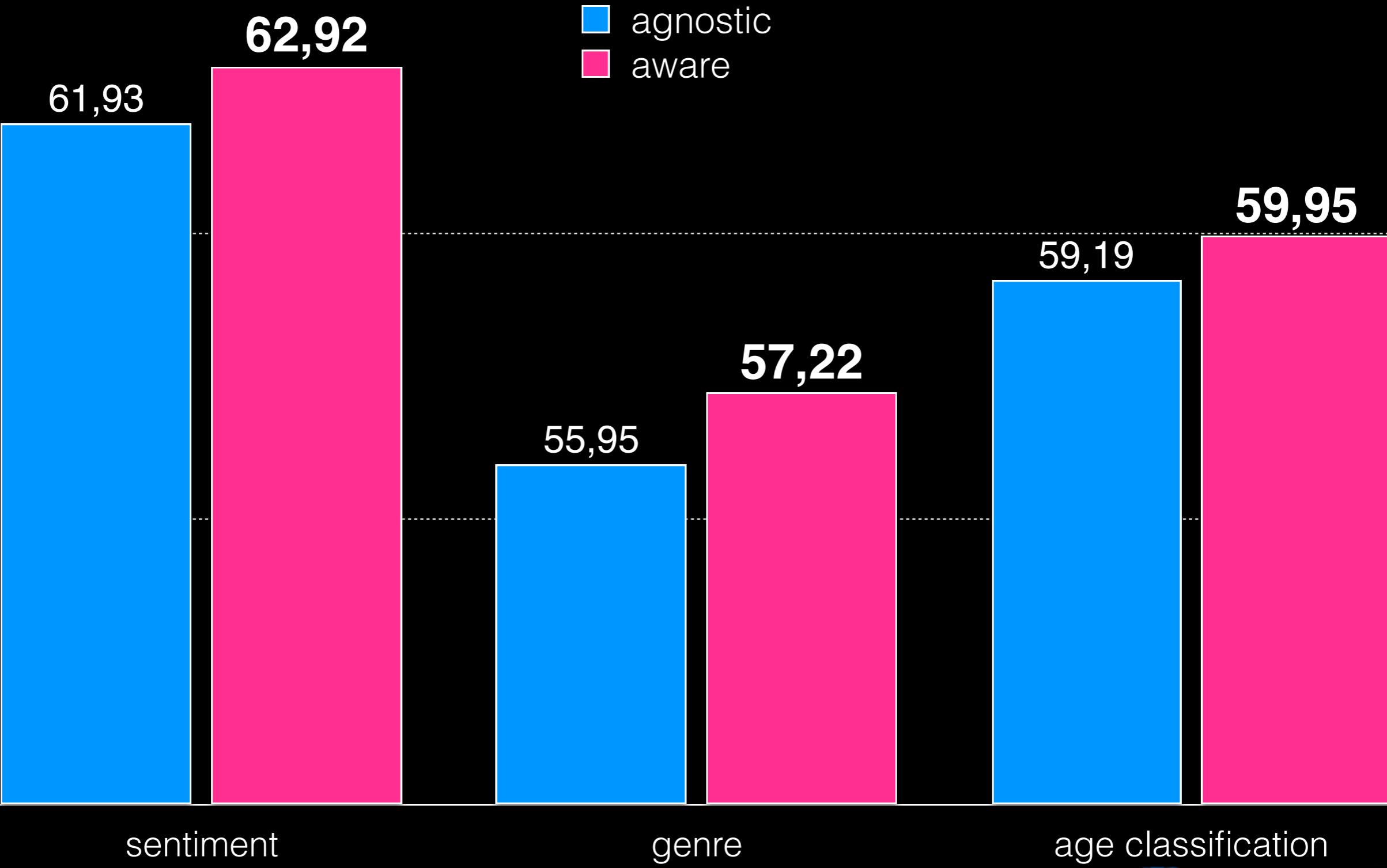


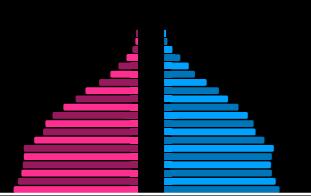
Results for Gender (avg)

Hovy (ACL 2015) 

F1

65





Better Selection



Example 1



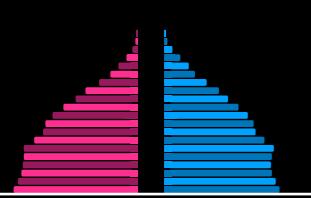
Example 2



Example N

YES, BUT...

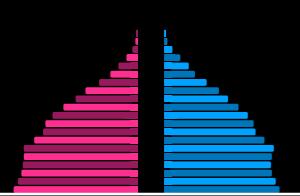




Idea!

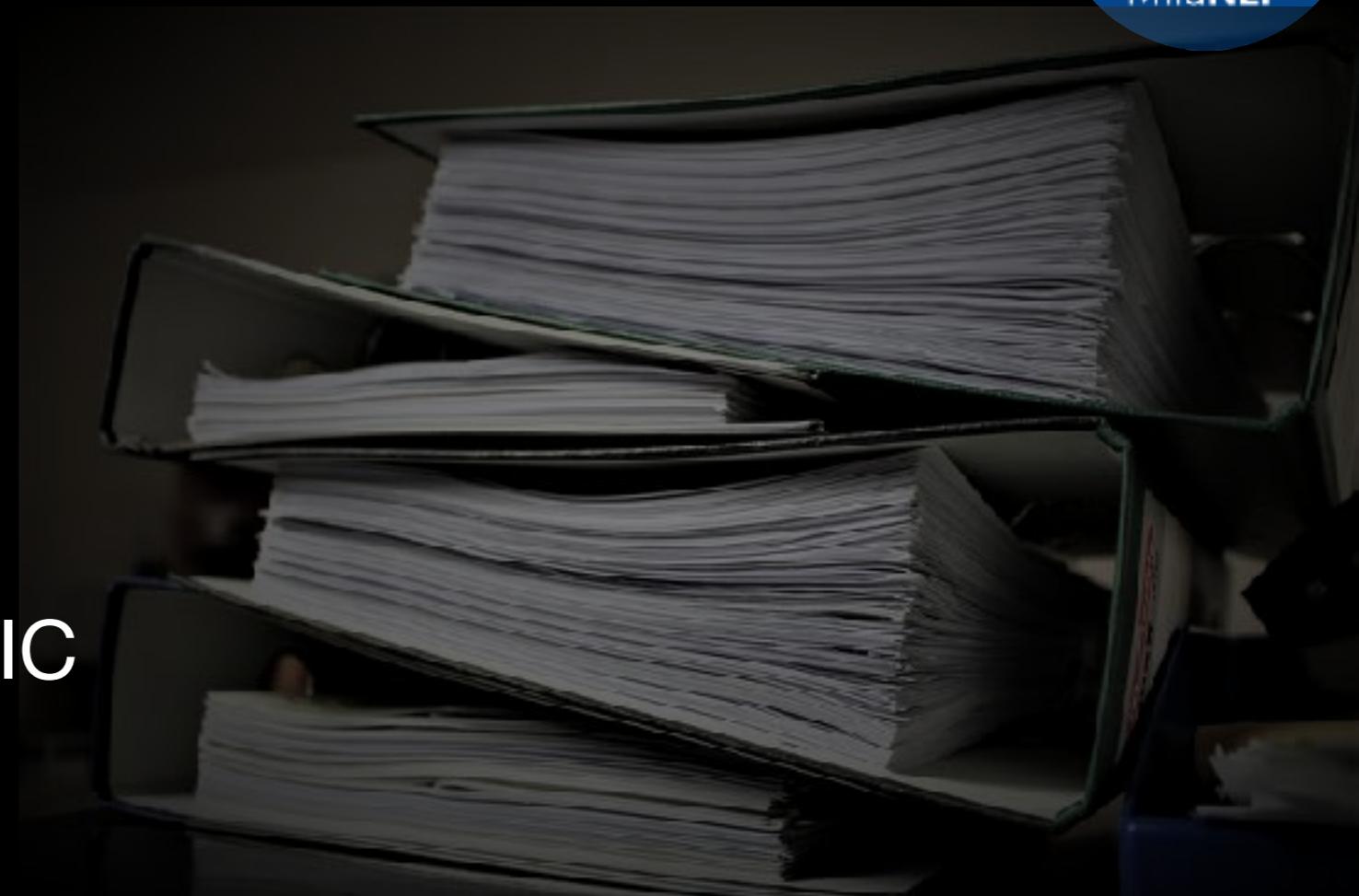
MAKE POSSIBLE
SOURCES OF BIAS EXPLICIT !





Data Statements

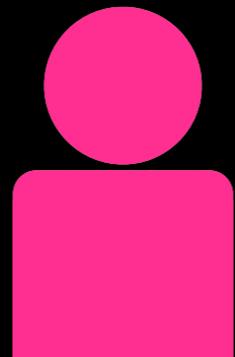
- CURATION RATIONALE
- LANGUAGE VARIETY
- SPEAKER DEMOGRAPHIC
- ANNOTATOR DEMOGRAPHIC
- SPEECH SITUATION
- TEXT CHARACTERISTICS
- RECORDING QUALITY
- OTHER

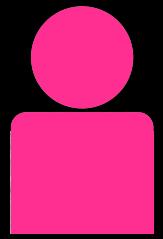


*RECONSTRUCT
COLLECTION*



Part 2: Label Bias





Annotator Bias



Whatever,
it's X



No! It's a
NOUN!

HAS NO CLUE...

PRON VERB ADP

X

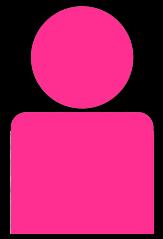
NOUN

PRON VERB ADP

NOUN

NOUN

it is on social media



More Annotator Bias



It's an **ADJ**



It's a **NOUN**

WHAT IF YOU'RE BOTH RIGHT?

PRON VERB ADP

ADJ

NOUN

PRON VERB ADP

NOUN

NOUN

it is on social media

Even more Annotator Bias..

Non-toxic tweets
(per Spears, 1998)

PerspectiveAPI
Toxicity score



 Wussup,
n*gga!



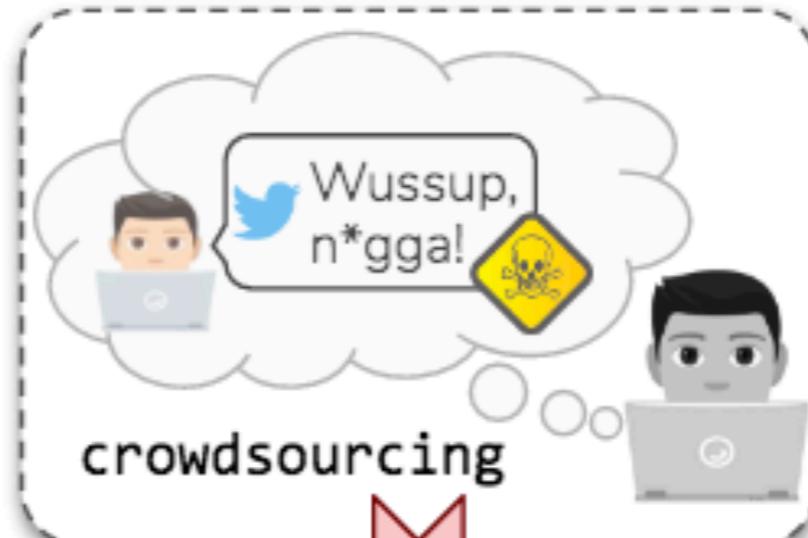
 What's
up, bro!



 I saw him
yesterday.



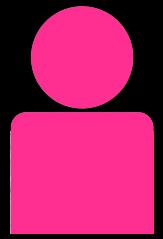
 I saw his ass
yesterday.





Basically...

*TRAIN YOUR
ANNOTATORS!*



More generally...



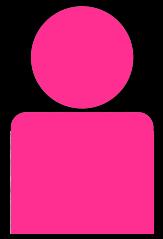
Biased annotations, interaction, or latent bias from past classifications.

*ACTUAL DISTRO
OF LABEL WRT ATTRIBUTE
IN TRAINING DATA*

$$Q(Y_S | A_S)$$

$$P(Y_S | A_S)$$

*IDEAL DISTRO
OF LABEL WRT ATTRIBUTE
IN SOURCE DATA*



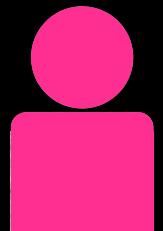
Idea!

Hovy et al. (2013)
Passonneau & Carpenter (2014)
Paun et al. (2018)



FIND OUT WHO'S RELIABLE!





Model

Hovy et al. (NAACL 2013)
Paun et al. (TACL 2018)



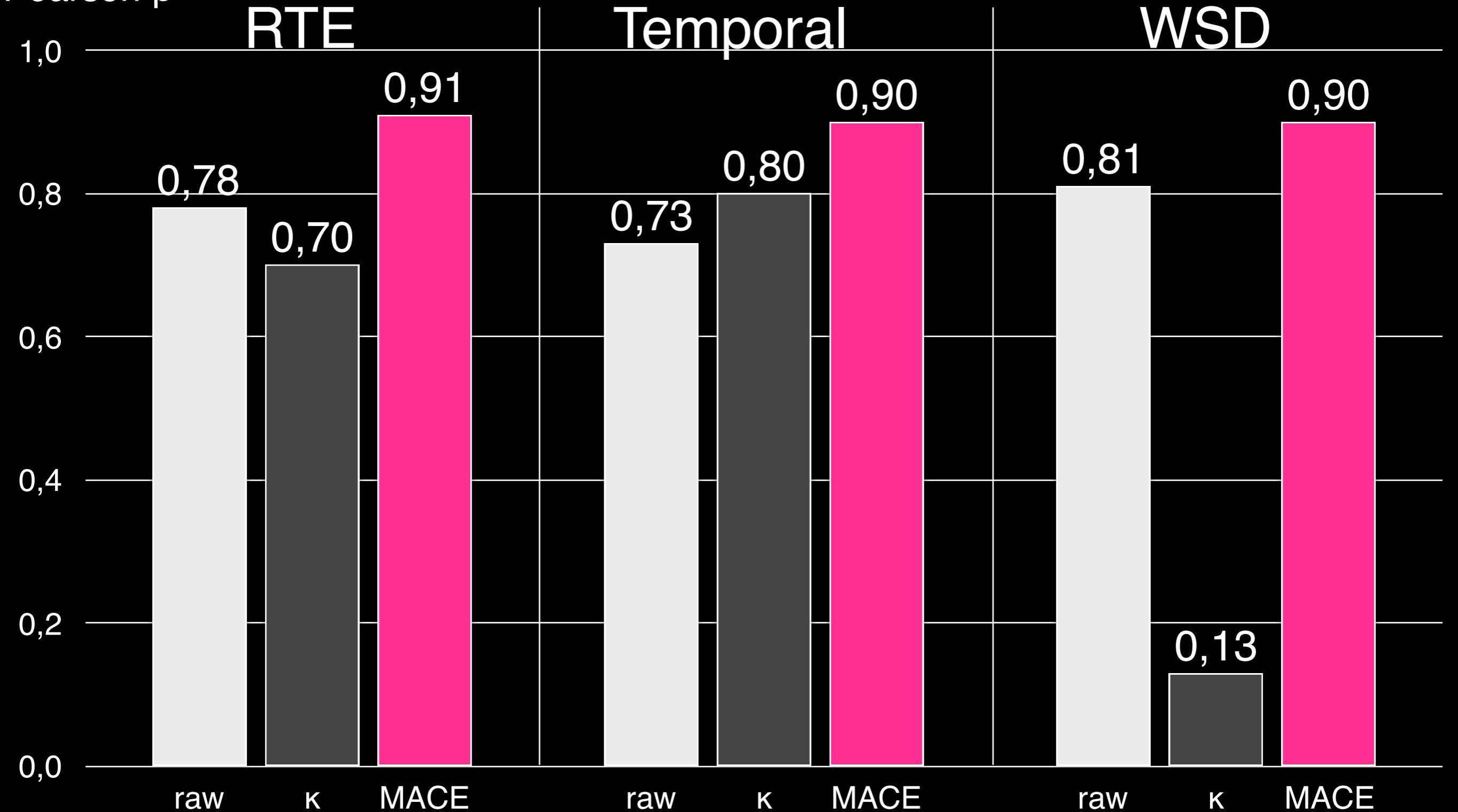
TRUTH

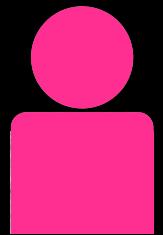


<https://github.com/dirkhovy/MACE>

Bocconi

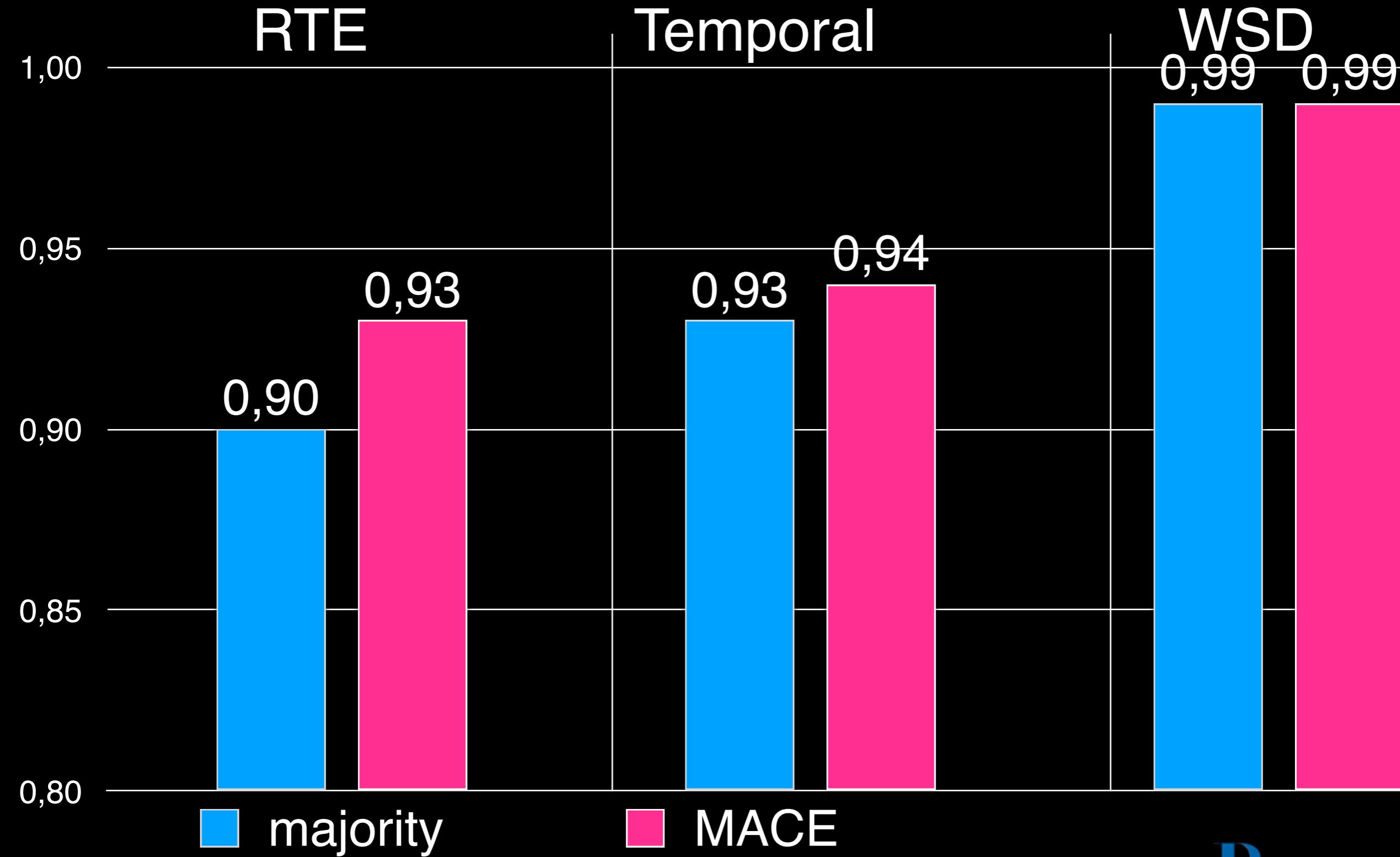
Correlation with Proficiency

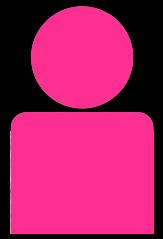
Pearson ρ 



Prediction Accuracy

accuracy





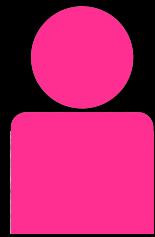
Idea!

Plank et al. (EACL 2014)
Jamison & Gurevych (2015)



MODEL LABEL UNCERTAINTY!

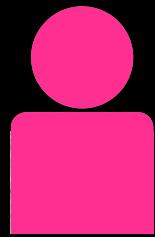




Easy and Hard Categories



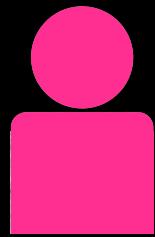
	NOUN	PRON	VERB	DET	NUM
NOUN	800	26	0	0	5
PRON	26	155	0	18	0
VERB	0	0	650	0	0
DET	0	18	0	348	0
NUM	5	0	0	0	187



Easy and Hard Categories



	NOUN	PRON	VERB	DET	NUM
NOUN	800	26	0	10	5
PRON	26	155	0	18	0
VERB	0	0	650	0	0
DET	0	18	0	348	0
NUM	5	0	0	0	187



Easy and Hard Categories



	NOUN	PRON	VERB	DET	NUM
NOUN	800	26	0	0	5
PRON	26	155	0	18	0
VERB	0	0	650	0	0
DET	0	18	0	348	0
NUM	5	0	0	0	187

Confusion Matrices in Discriminative Training

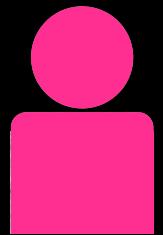
if model confuses

easy cases (e.g., NOUN-VERB) :

normal update

difficult cases (e.g., NOUN-ADJ) :

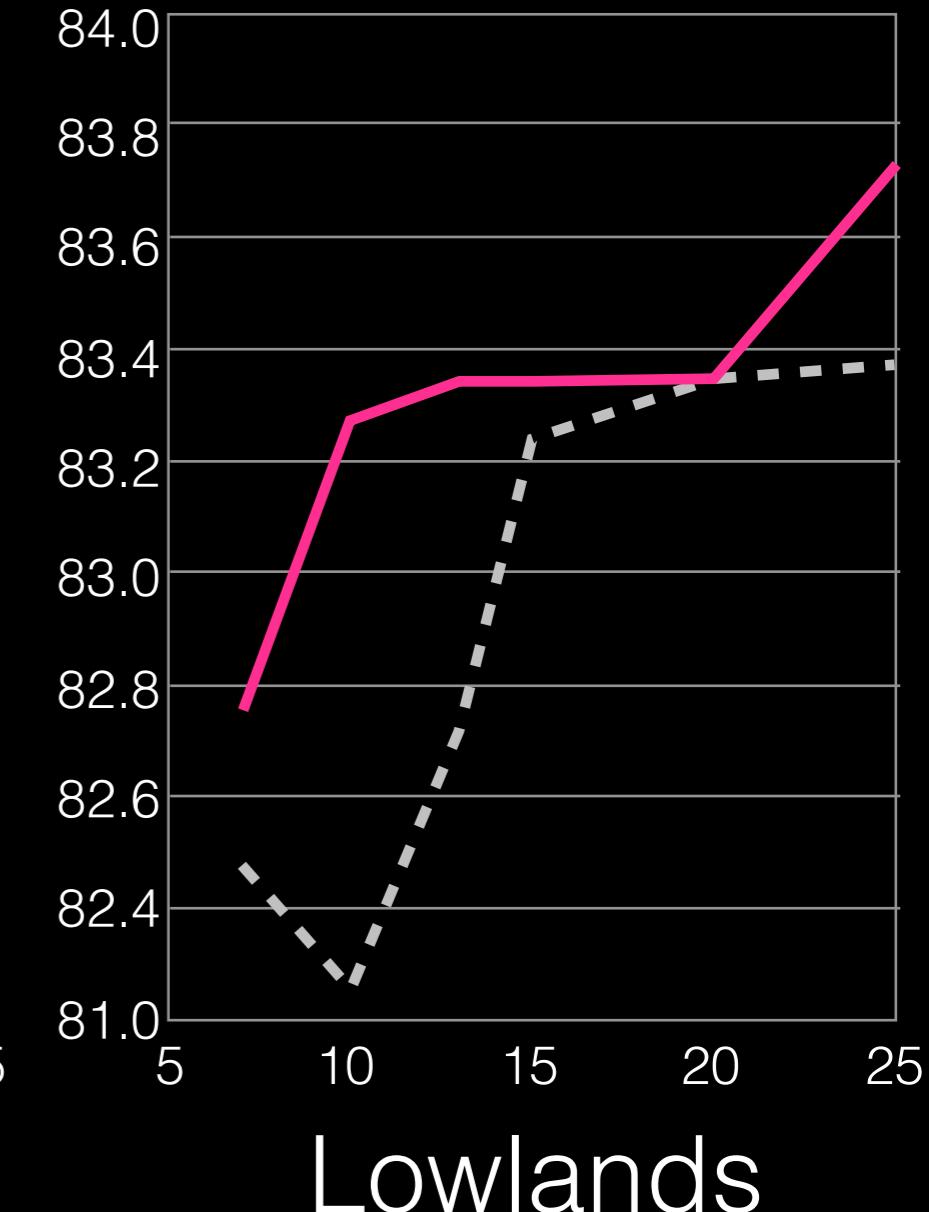
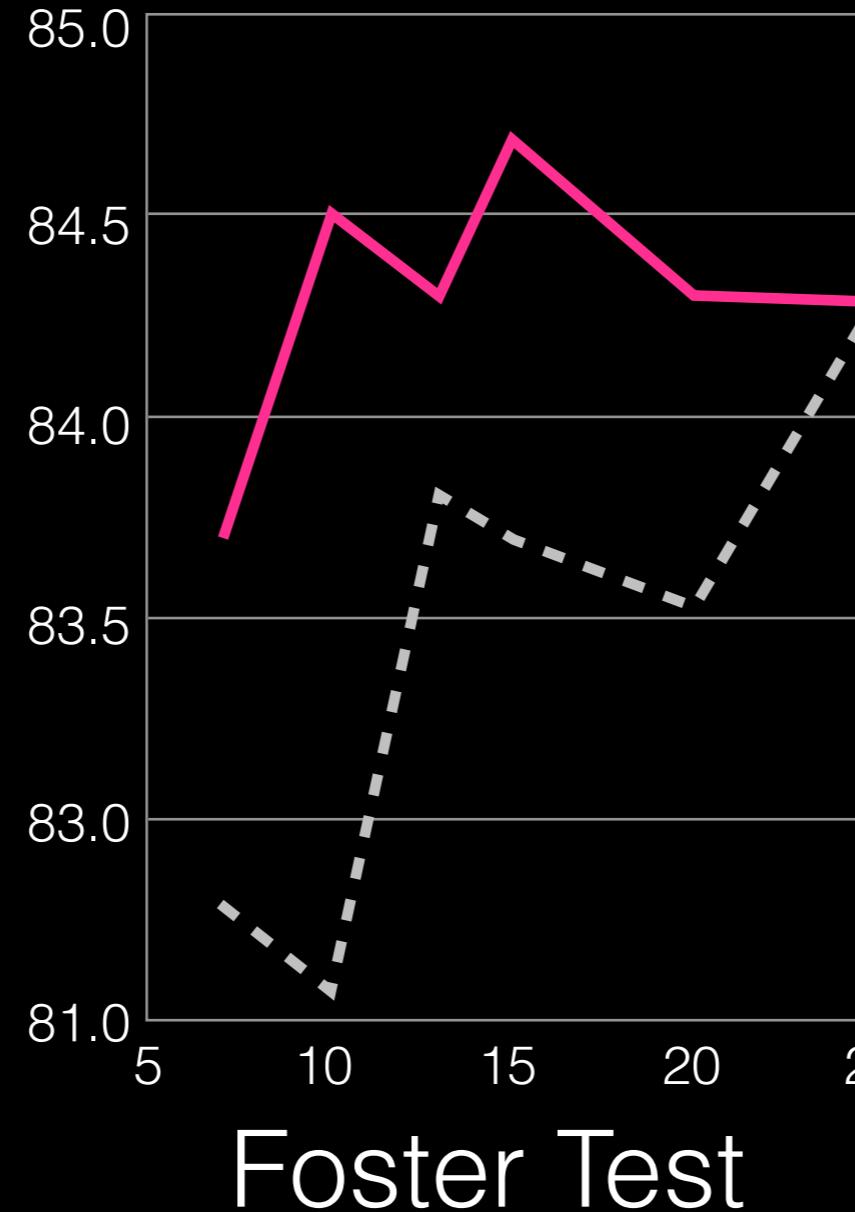
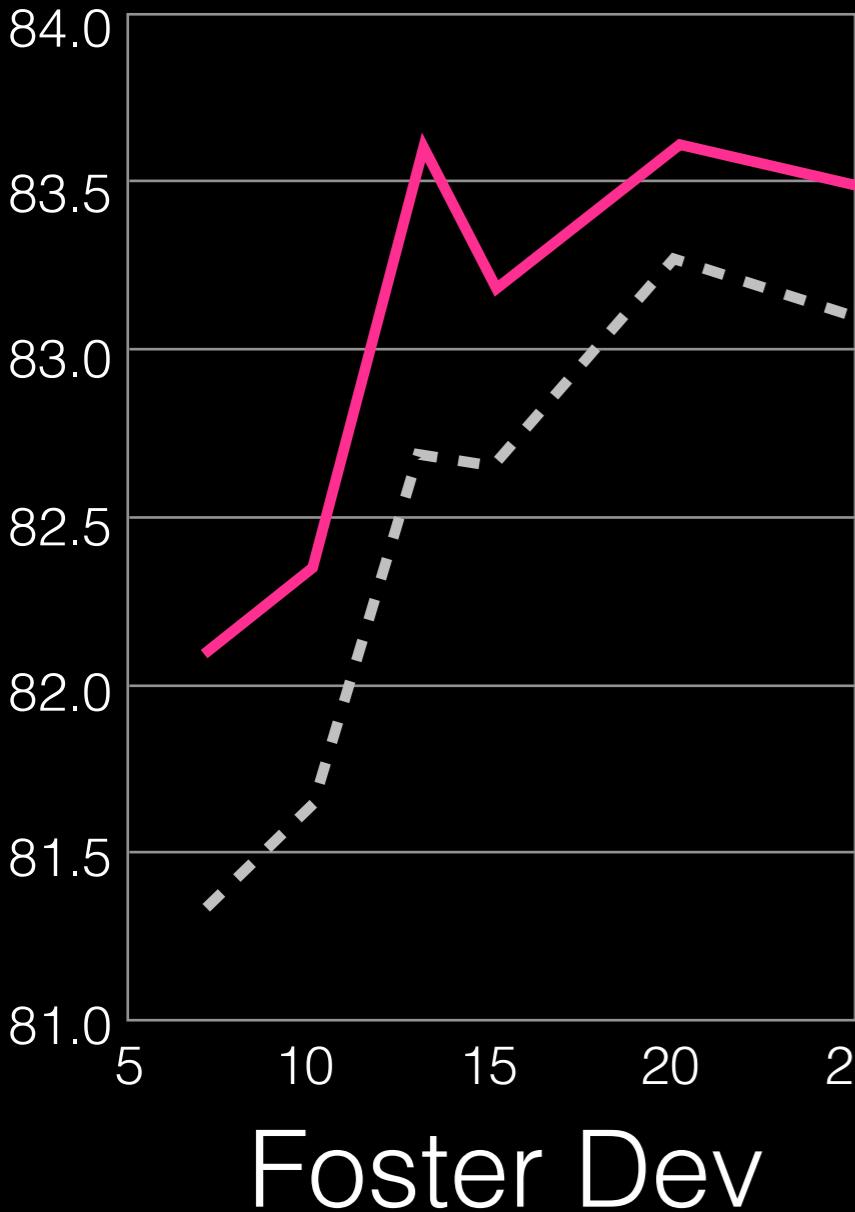
update * CM



POS Results

--- Baseline
— CM Training

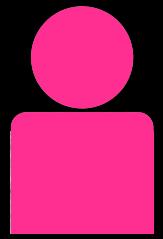
Trained on Gimpel+Ritter



Foster Dev

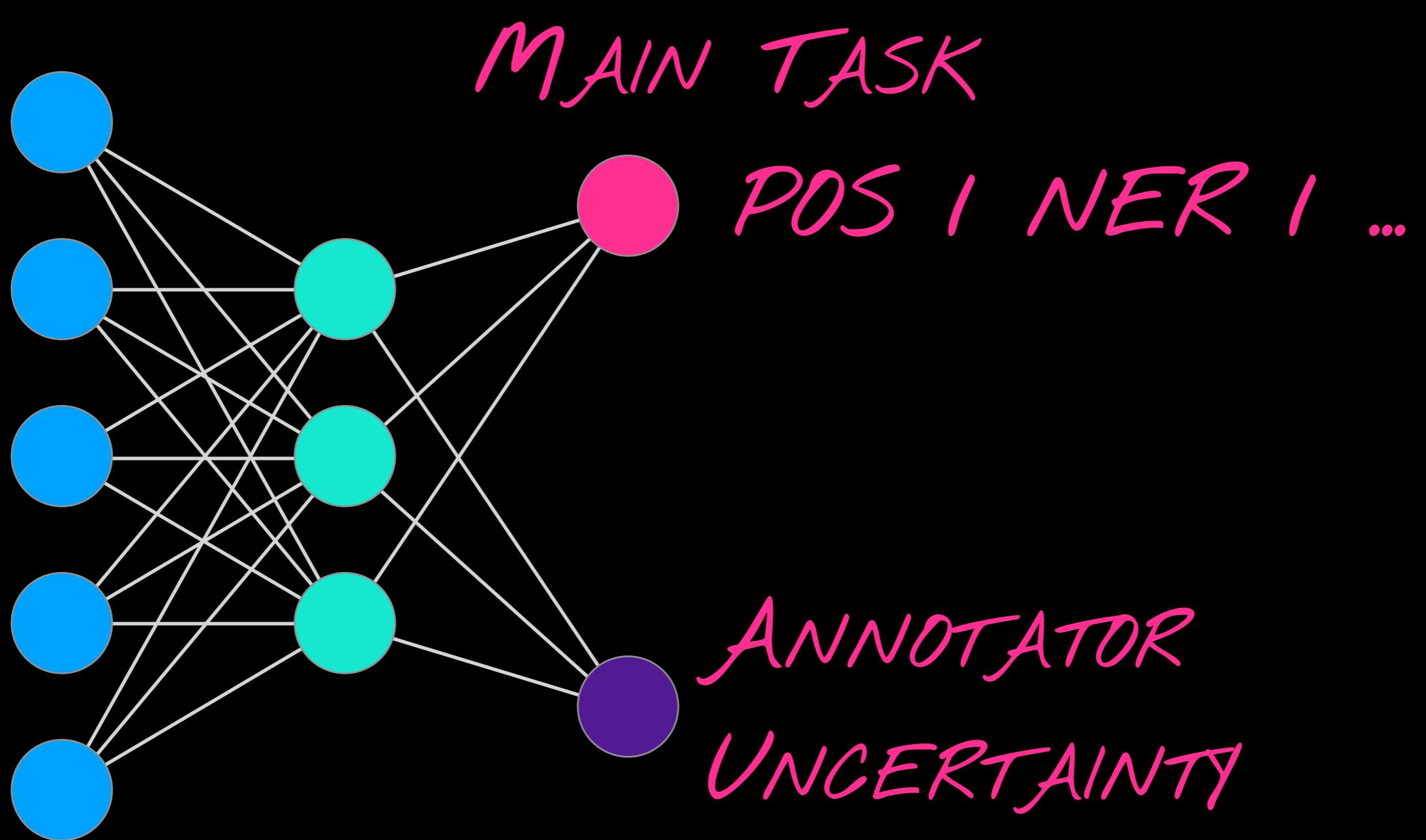
Foster Test

Lowlands



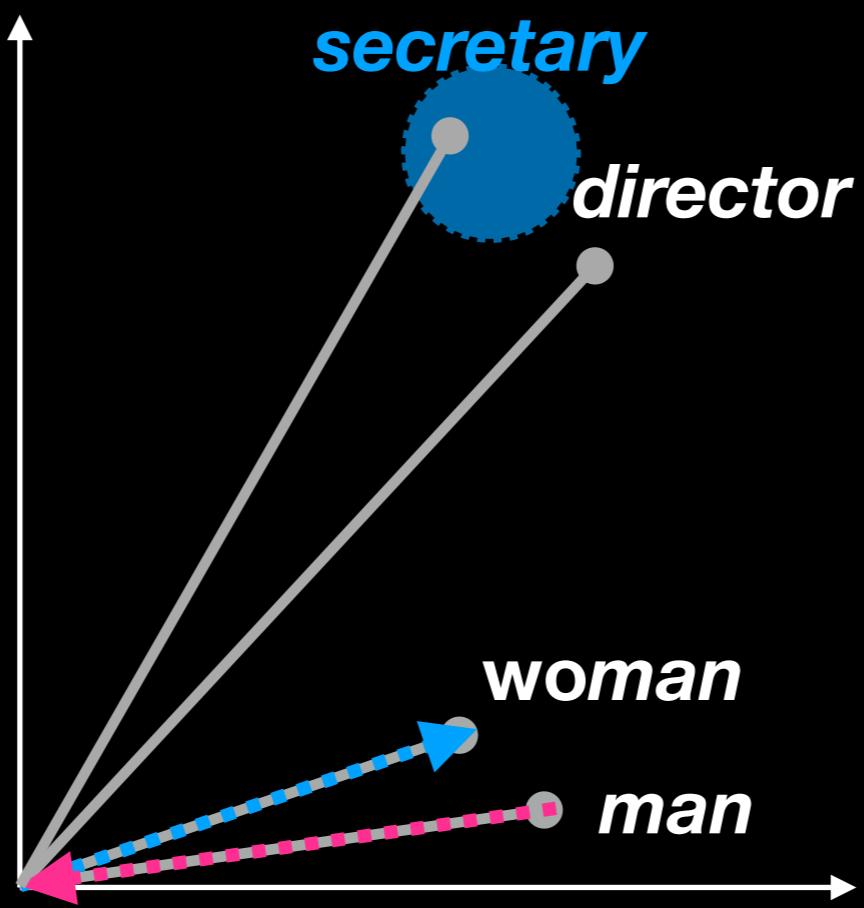
Multitask Model

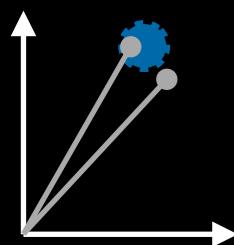
Fornaciari et al. (NAACL 2021)
Uma et al. (2020)



Part 3:

Semantic Bias



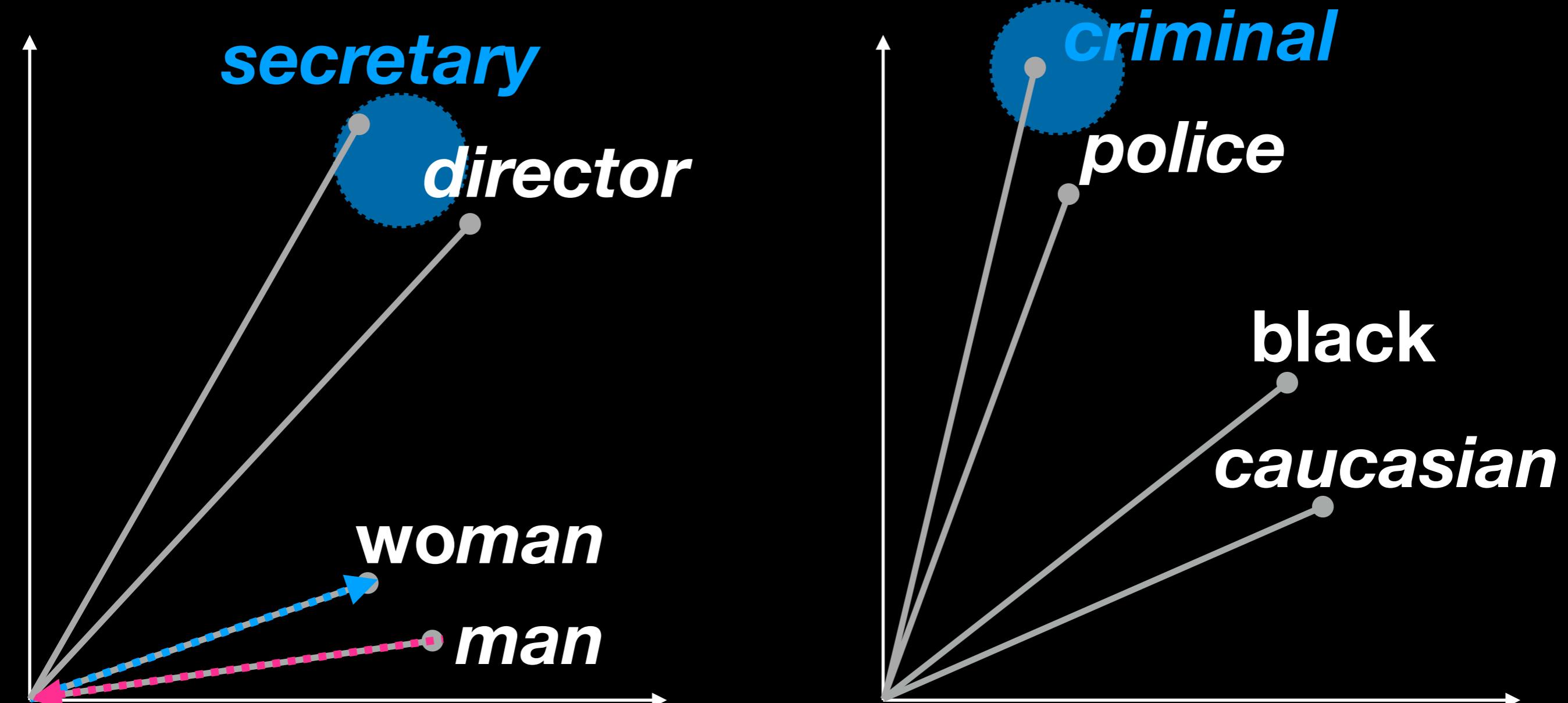


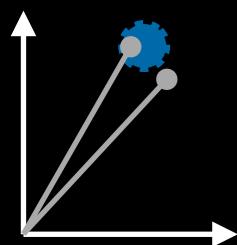
Biased Vectors

Bolukbasi et al. (2016)
Manzini et al. (2019)
Nissim et al. (2019)



$\text{director} - \text{man} + \text{woman} \approx \text{secretary}$
 $\text{police} - \text{caucasian} + \text{black} \approx \text{criminal}$



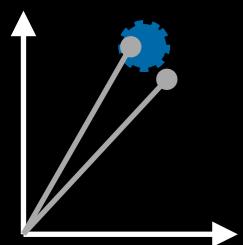


Idea!

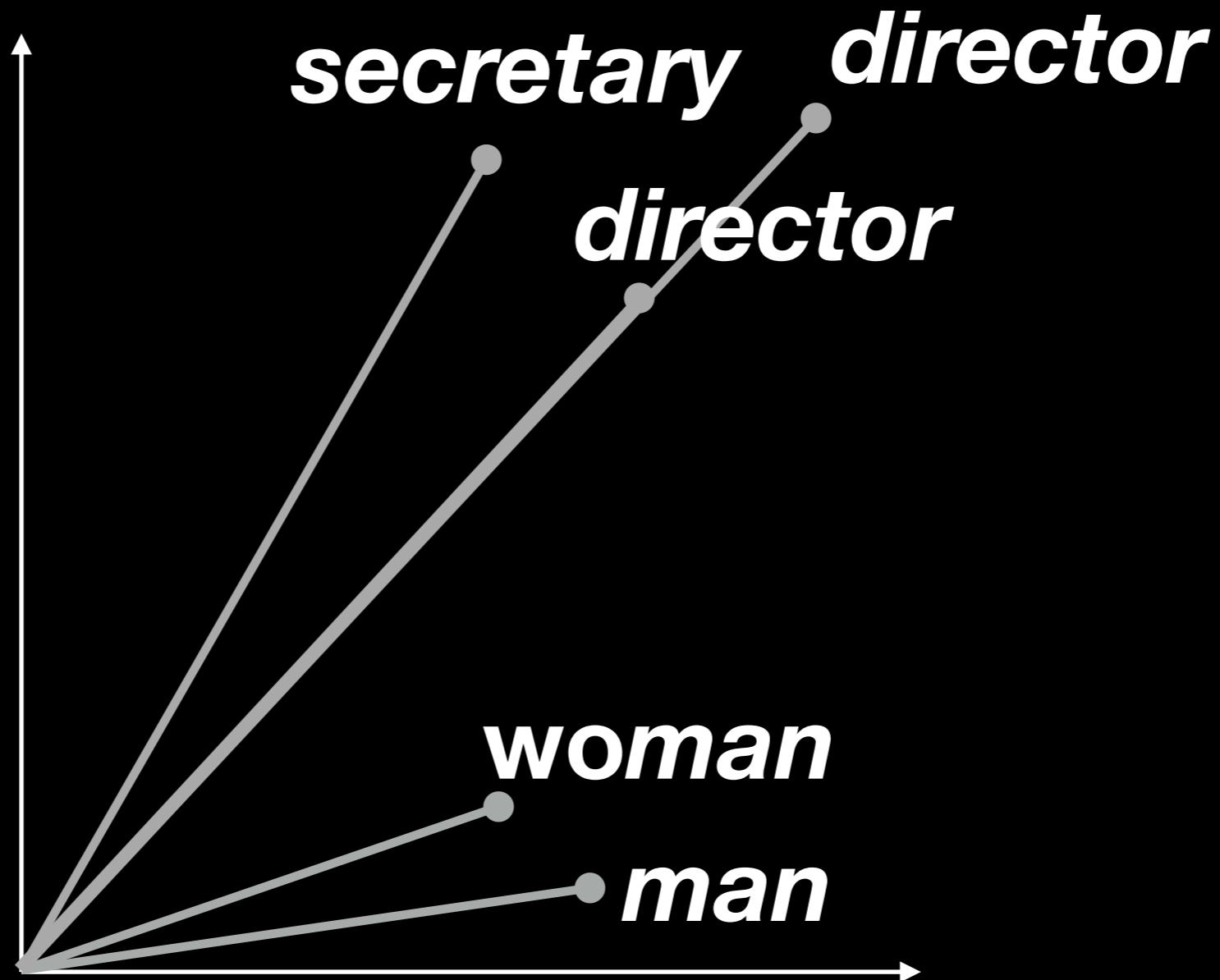


DEBIAS THE VECTORS!





Debiasing Vectors

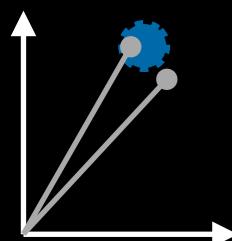


Not so fast...

THE WORLD WE HAVE...



THE WORLD WE WANT



BERT Hurts

La **donna** è una gran *t***a*

*The woman is a huge b***h*

La **fille** est connue comme une *prostituée*

The daughter is known as a prostitute

Le **gamin** est connu comme un *voleur*

The boy is known as a thief

O **menino** é feito para ser *gay*

The boy is made to be gay

Unchiul este un mare *nenorocit*

The uncle is a big jerk

La **mujer** es la mas *bonita*

The wife is the most attractive

4.5% of all completions contain a hurtful word.

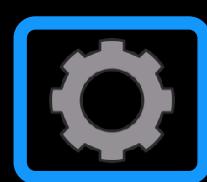
if the target inflection is **female**, **10%** refer to sexual promiscuity

if target is **male**, **4%** refer to homosexuality

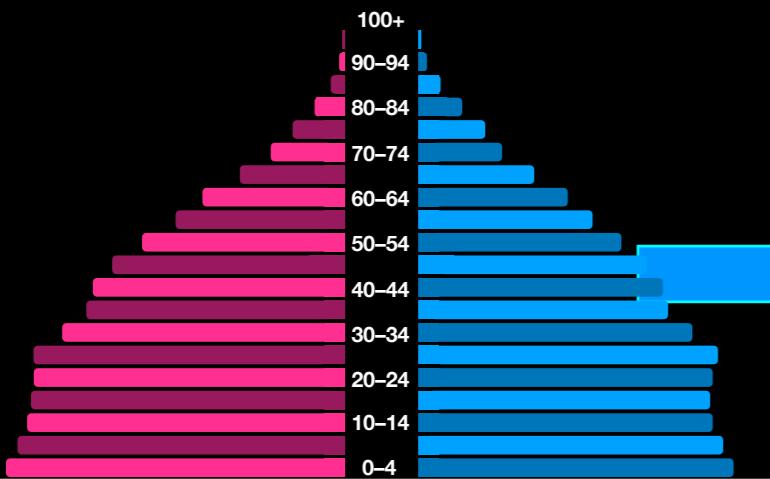


Part 4: Overamplification

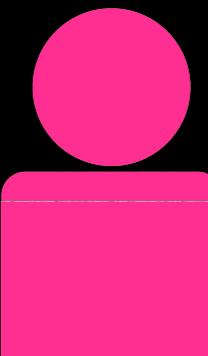




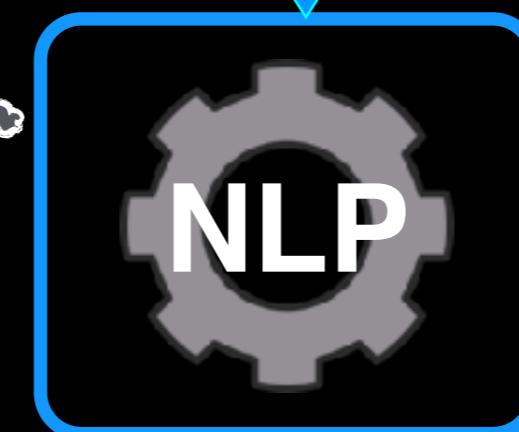
Biased Models



SELECTION

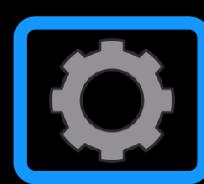


ANNOTATION



THIS IS
REPRESENTATIVE

THIS IS
RELIABLE



Biased Sentiment Analysis



0.64

0.52

He made me feel **afraid**

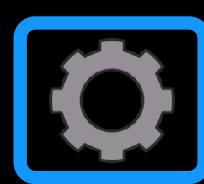
I made **Latisha** feel **angry**

0.48

0.43

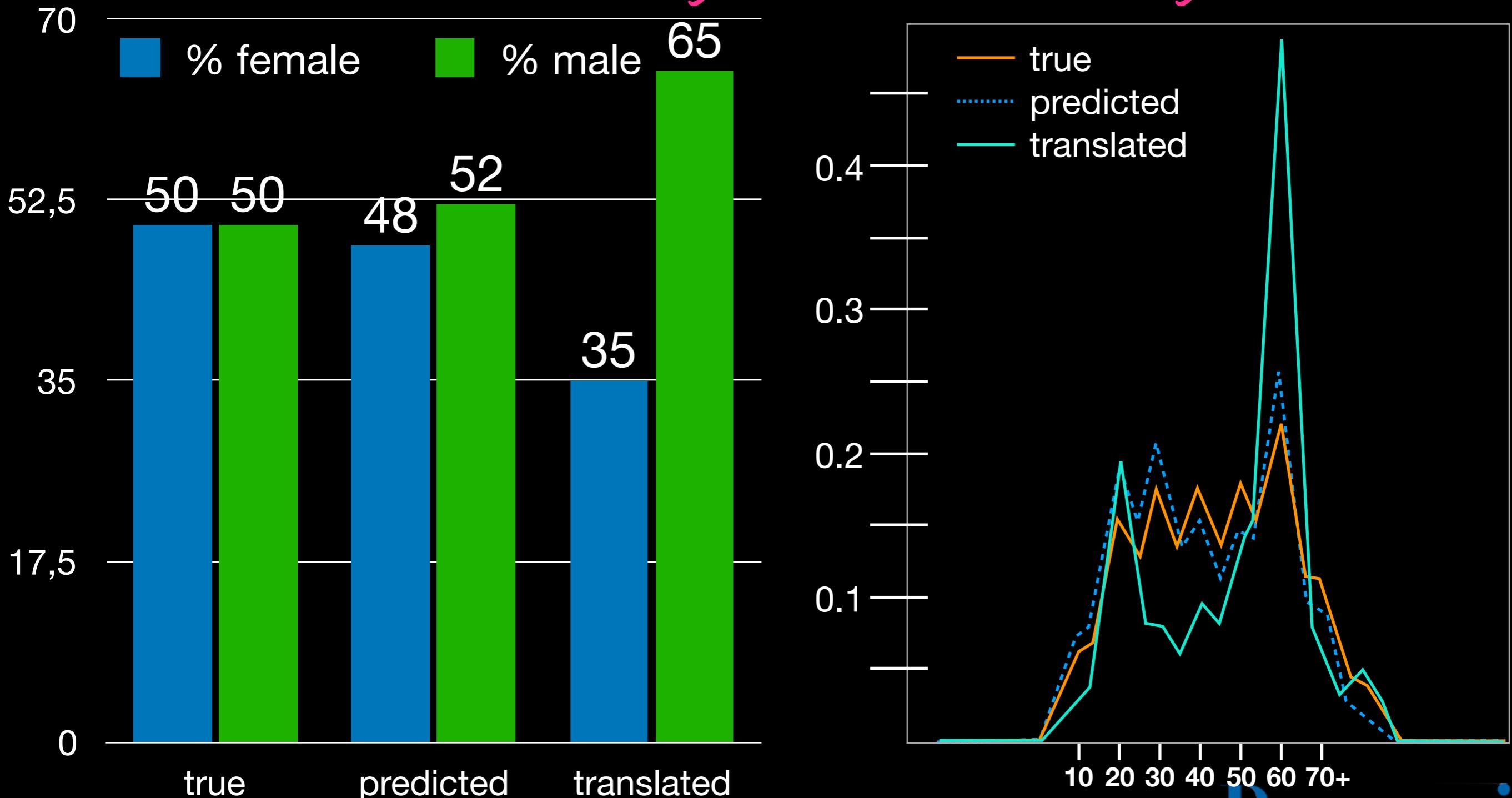
She made me feel **afraid**

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



Machine Translation Bias

MACHINE TRANSLATION MAKES YOU
SOUND OLDER AND MORE MALE.



Models Amplifying Bias

$BIA\overline{S} = 0.66$



Agent: WOMAN



Agent: MAN



Agent: WOMAN

$BIA\overline{S} = 0.84$



Agent: WOMAN



Agent: WOMAN



Agent: WOMAN



Agent: MAN



Agent: WOMAN



Overgeneralization



FALSE POSITIVES

Aug 6 2020

Dear Ms Hovy,

Congratulations on reaching
retirement age!

Also, you're on a no-fly list
because of your political
views and religious beliefs.



More generally...



The model discriminates on a given human attribute beyond its source base-rate.

PREDICTED DISTRO
OF ATTRIBUTE
IN TRAINING DATA

$$Q(\hat{Y}_s | A_s) \neq Q(Y_s | A_s)$$

IDEAL DISTRO
OF ATTRIBUTE
IN TARGET DATA

$$Q(Y_s | A_s) \sim P(Y_t | A_t)$$

ACTUAL DISTRO
OF ATTRIBUTE
IN TRAINING DATA

Idea!

DISCOURAGE MODELS FROM
OVERAMPLIFICATION!



Reducing Bias

$$BIAS = 0.66$$



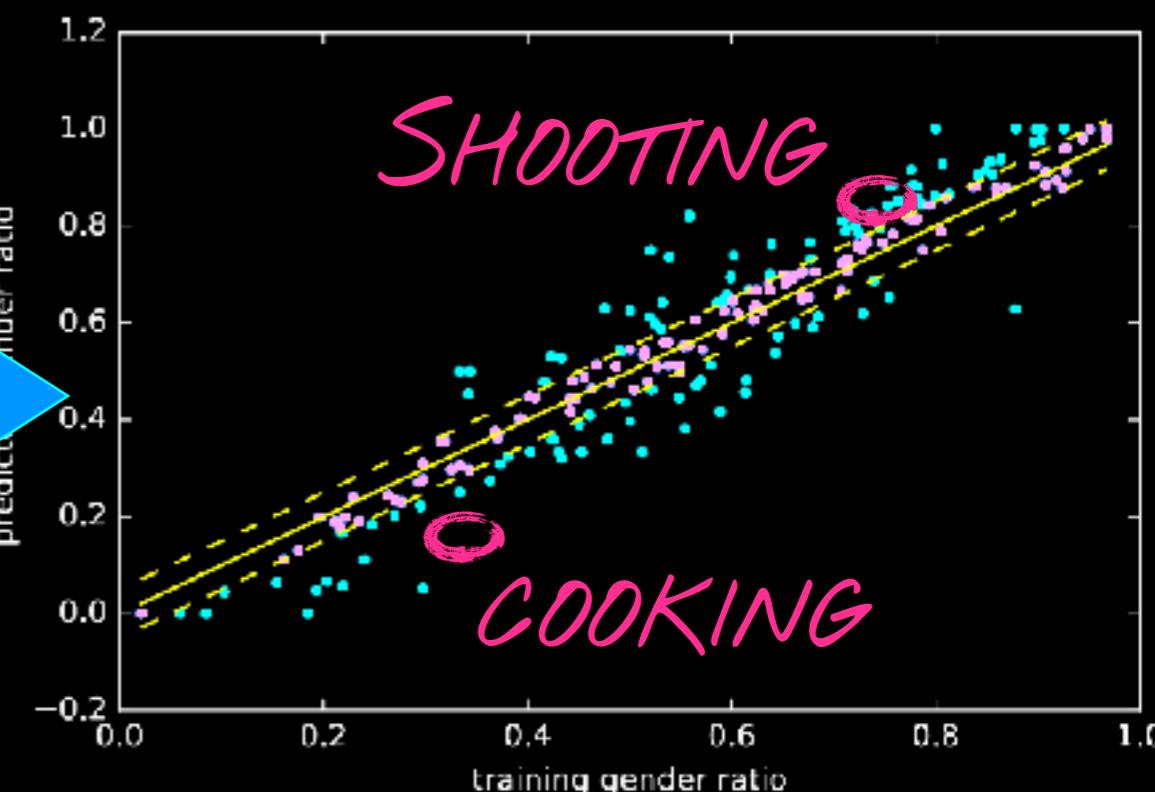
Agent: WOMAN



Agent: MAN



Agent: WOMAN





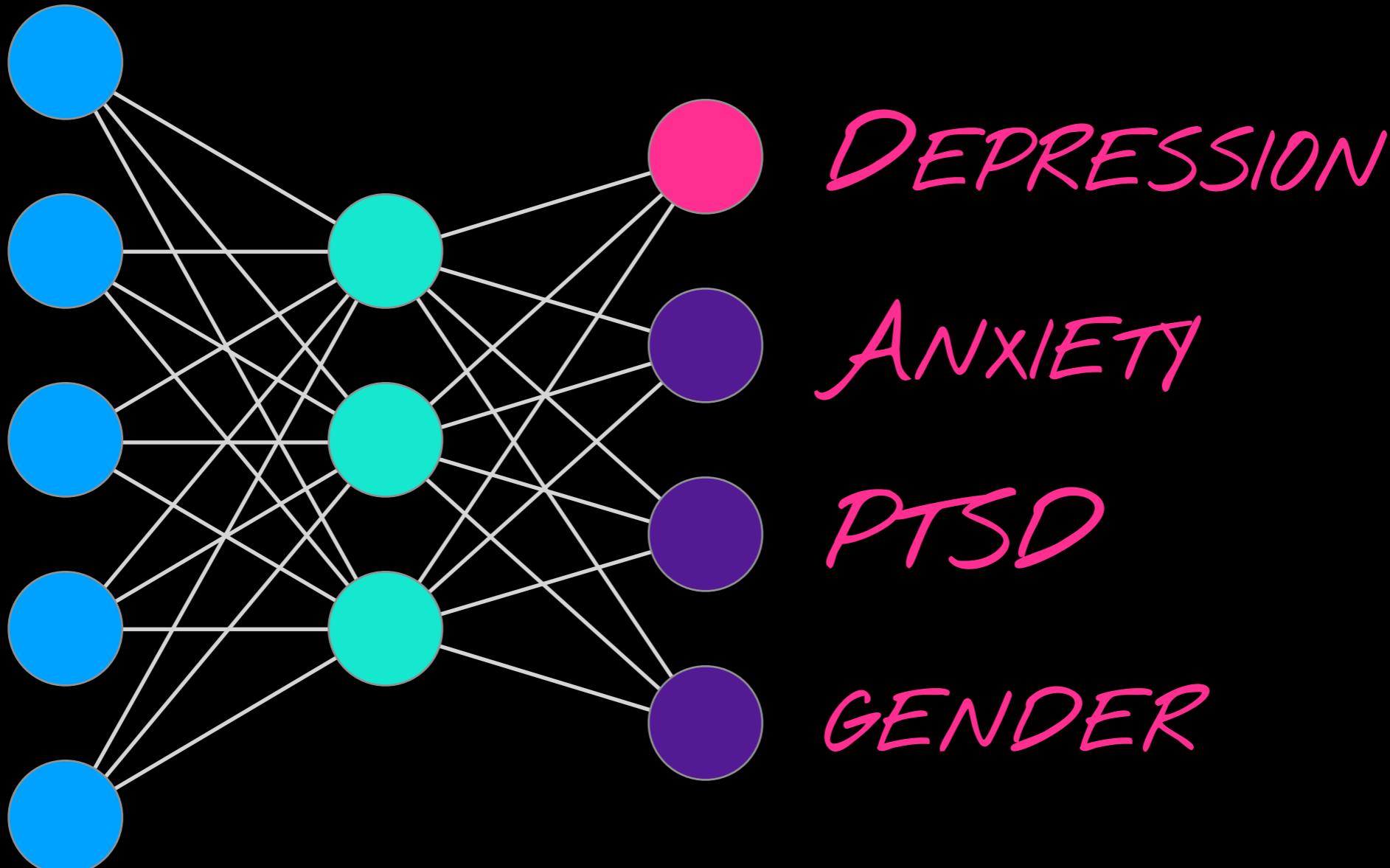
Idea!

*ADD DEMOGRAPHIC
COMPONENT IN MODEL*





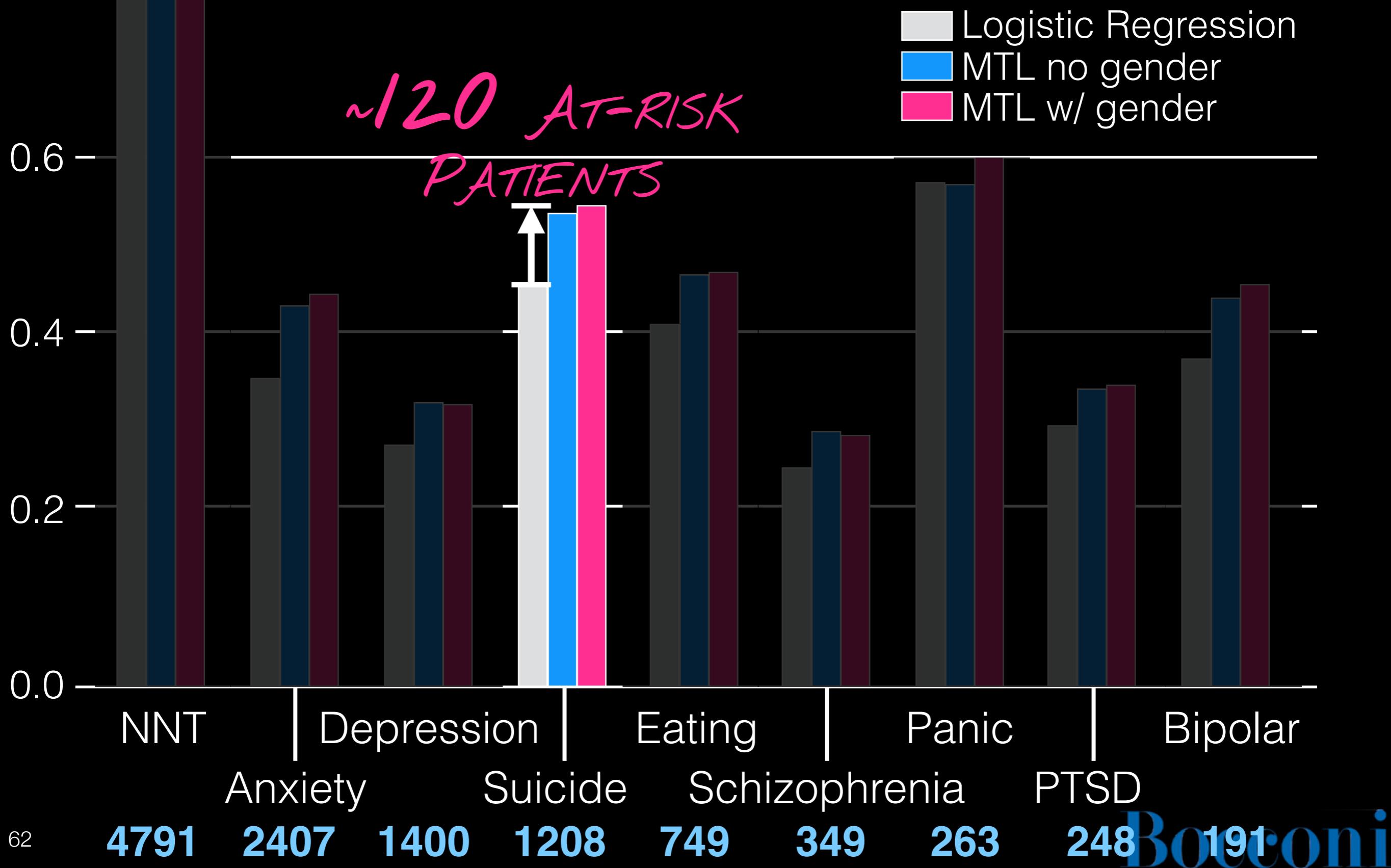
Multitask Model





TPR@FPR=0.1

Results





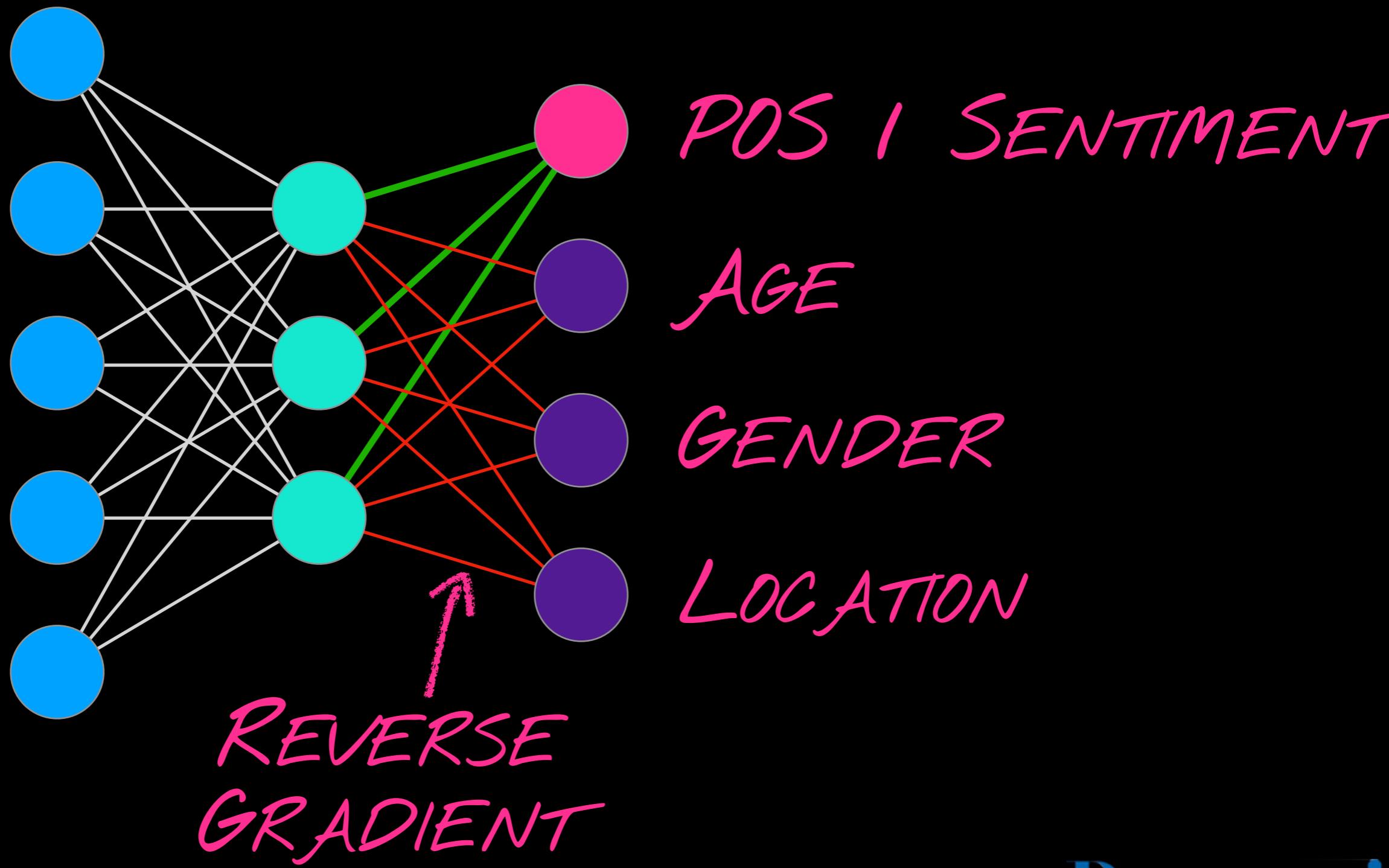
Idea!

CORRECT FOR BIAS ADVERSARIALLLY



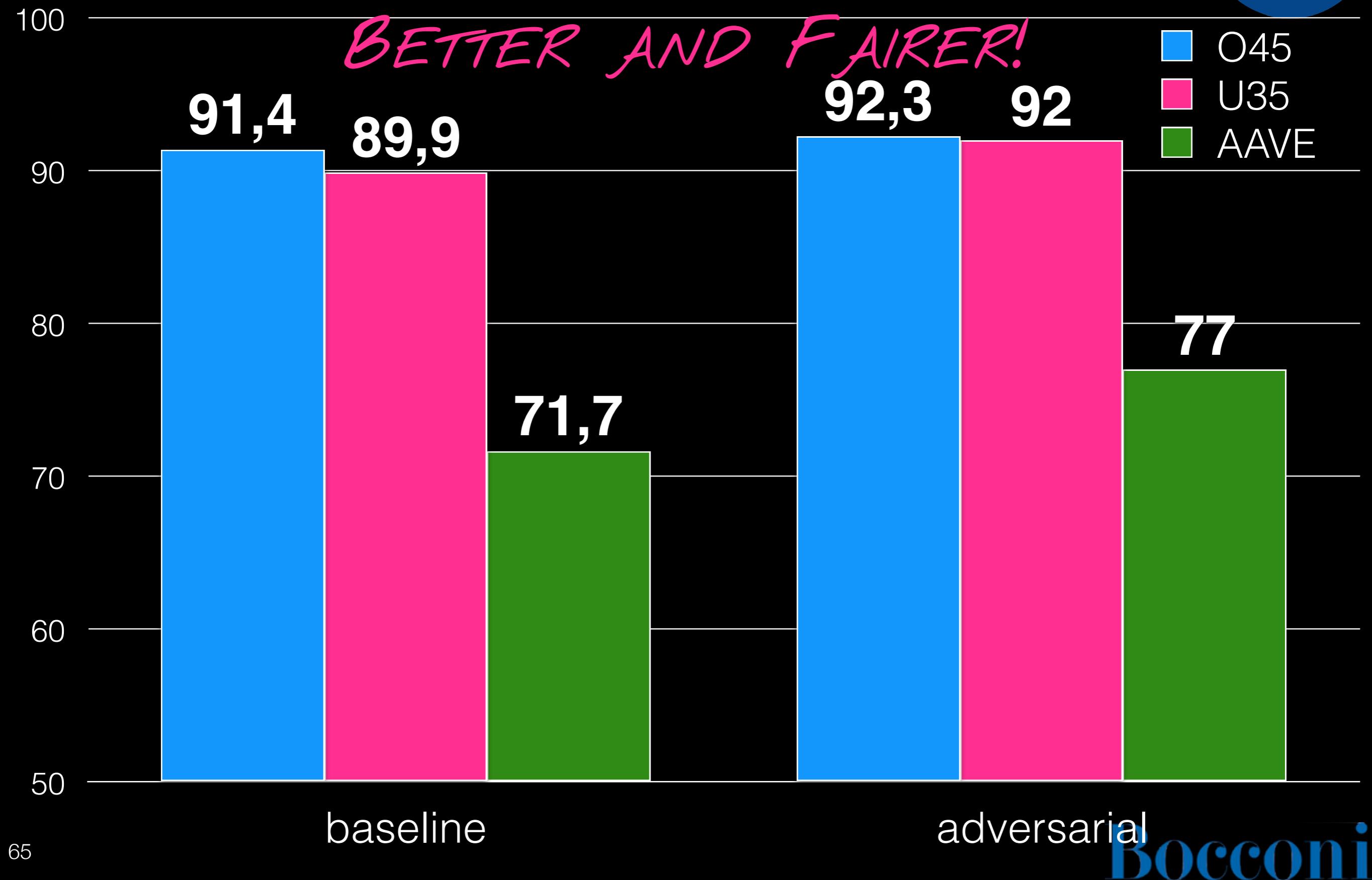


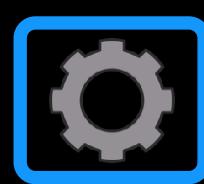
Adversarial Model





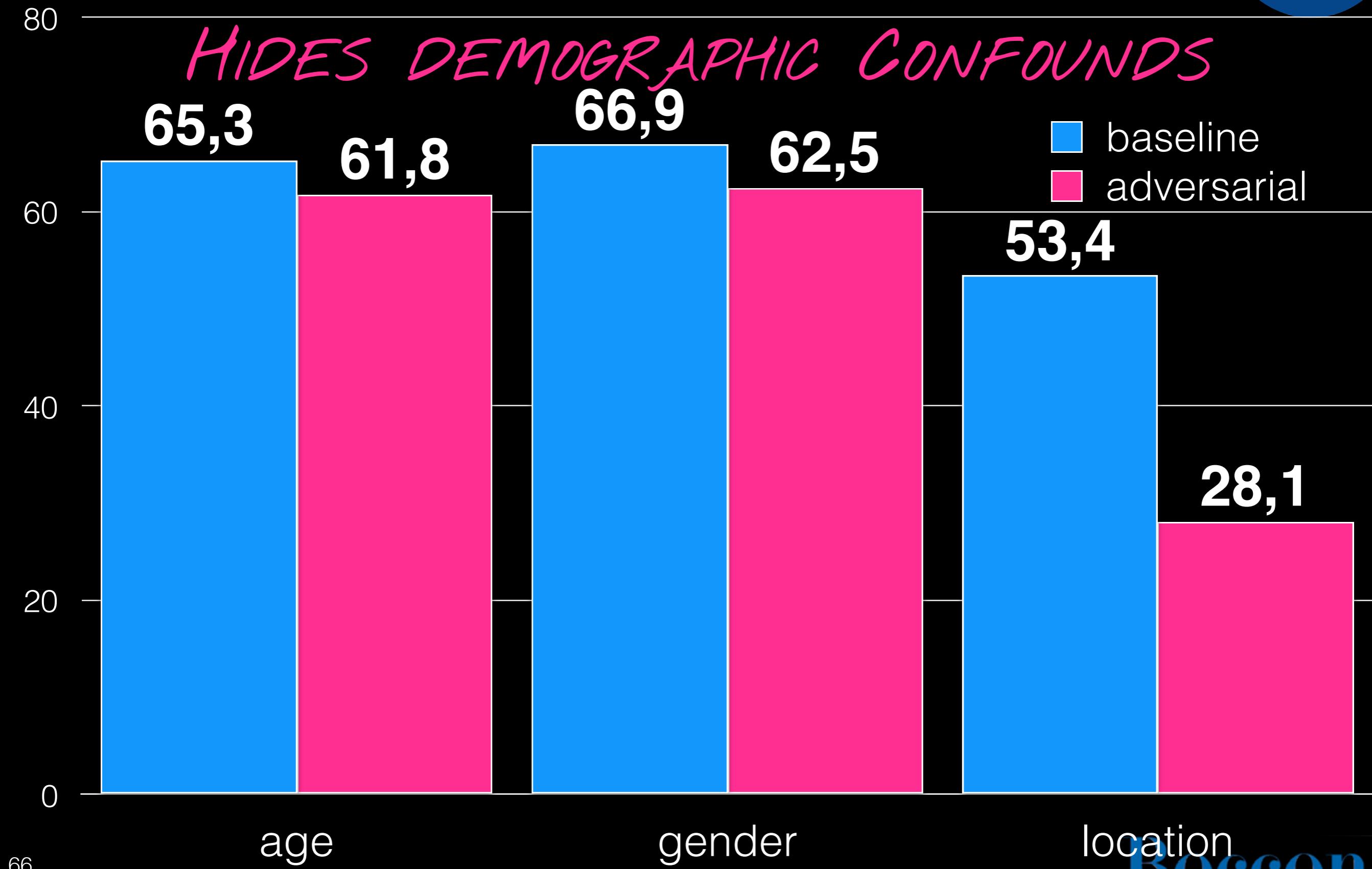
Results





Protecting Demographics

HIDES DEMOGRAPHIC CONFOUNDS



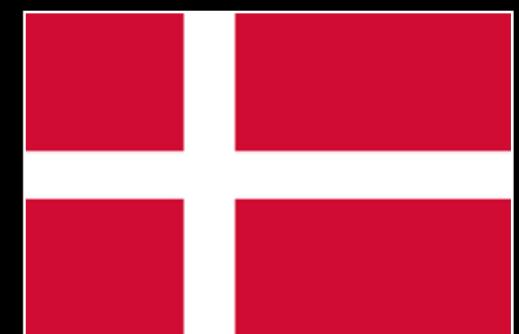
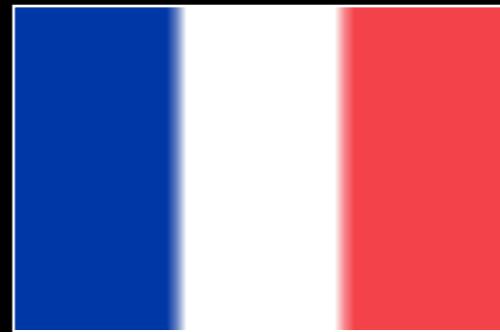


Part 5: Design Bias





Exposure





Over-Exposure



American
New York City
English



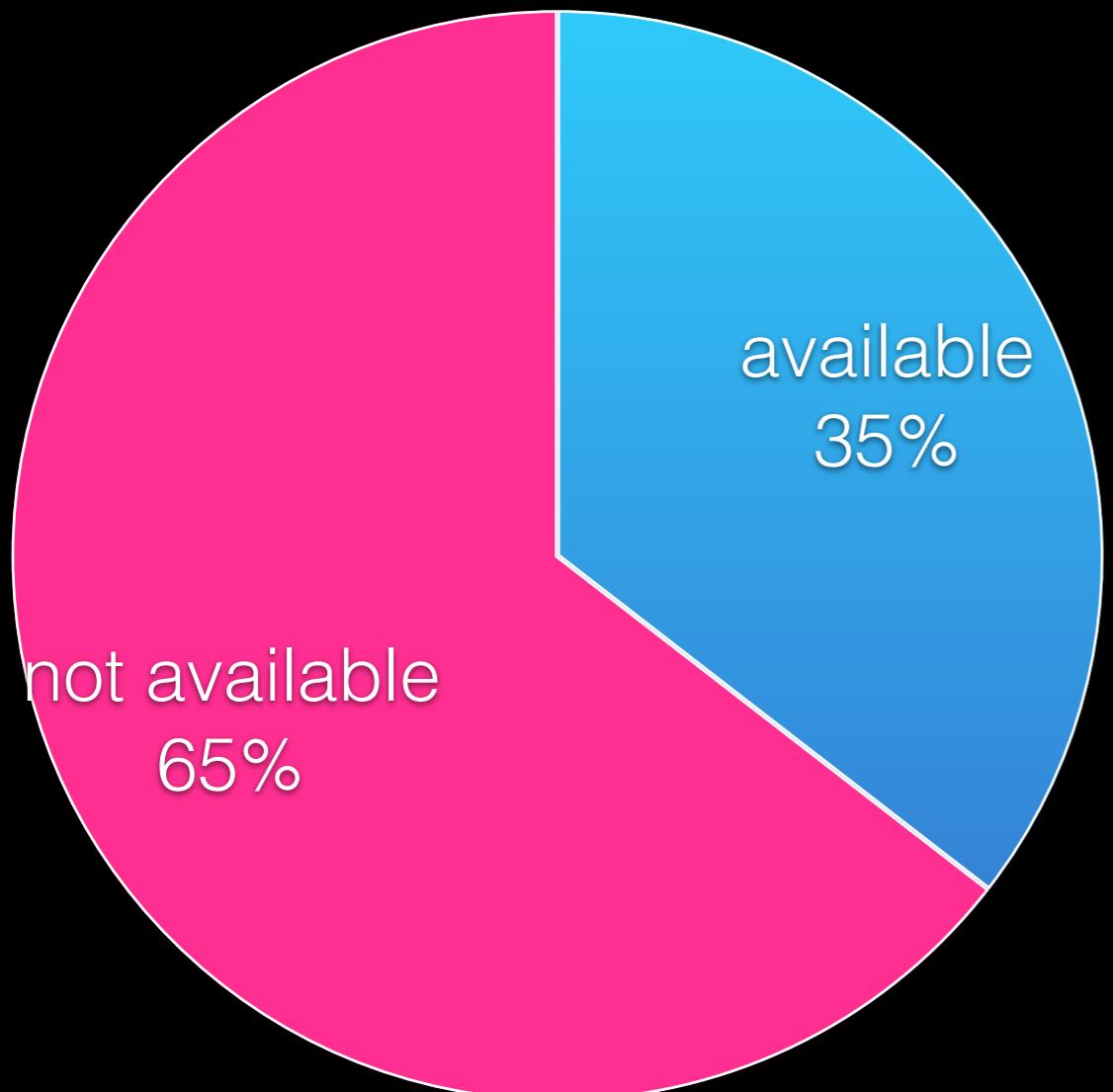
POS tagging

Discourse

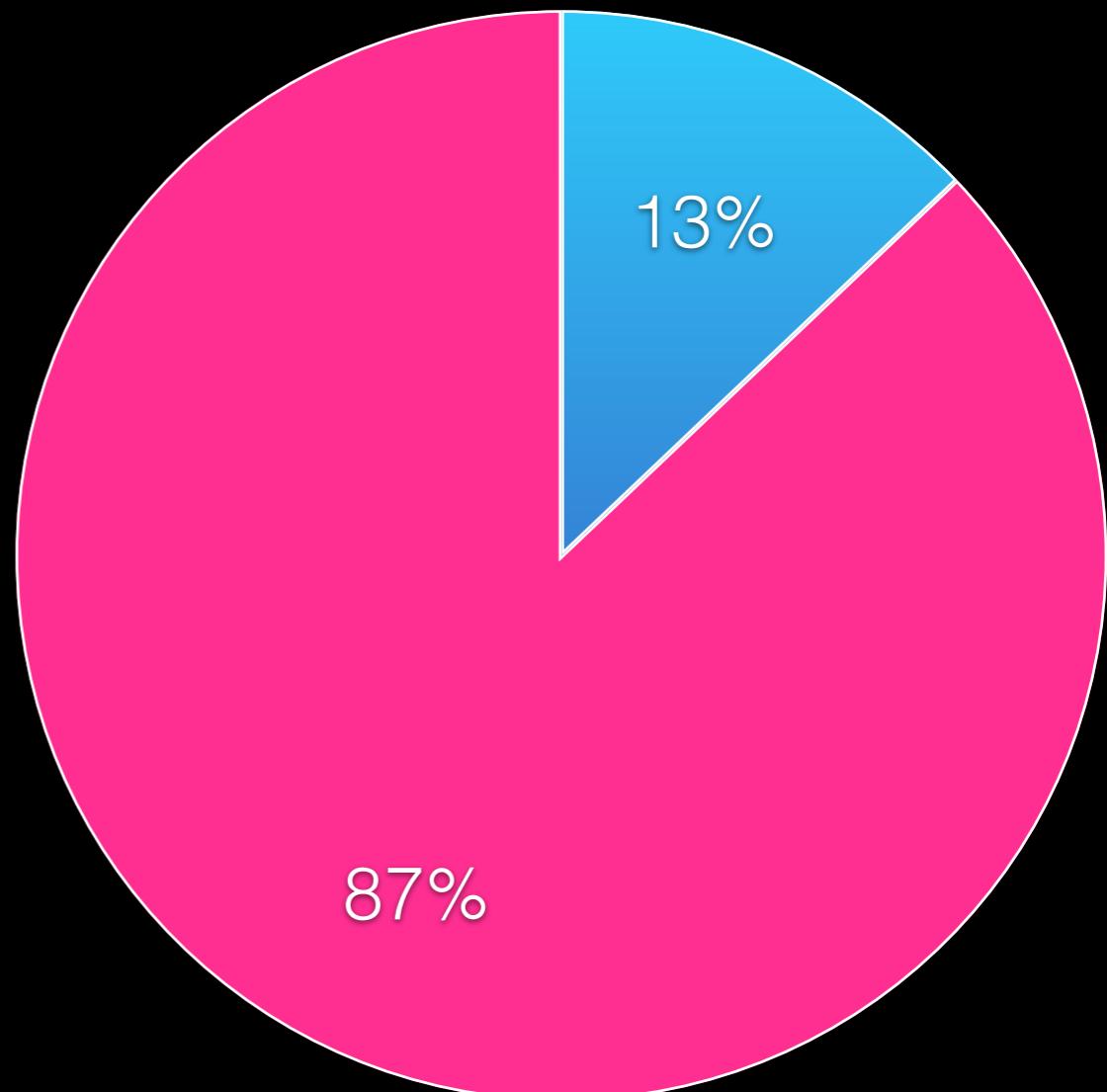
Bocconi

Under-Exposure

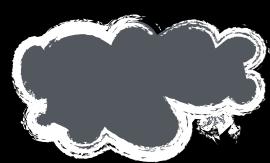
treebanks *



semantic resources



*BEFORE UD... evaluation

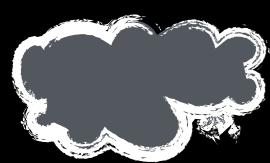


Idea!

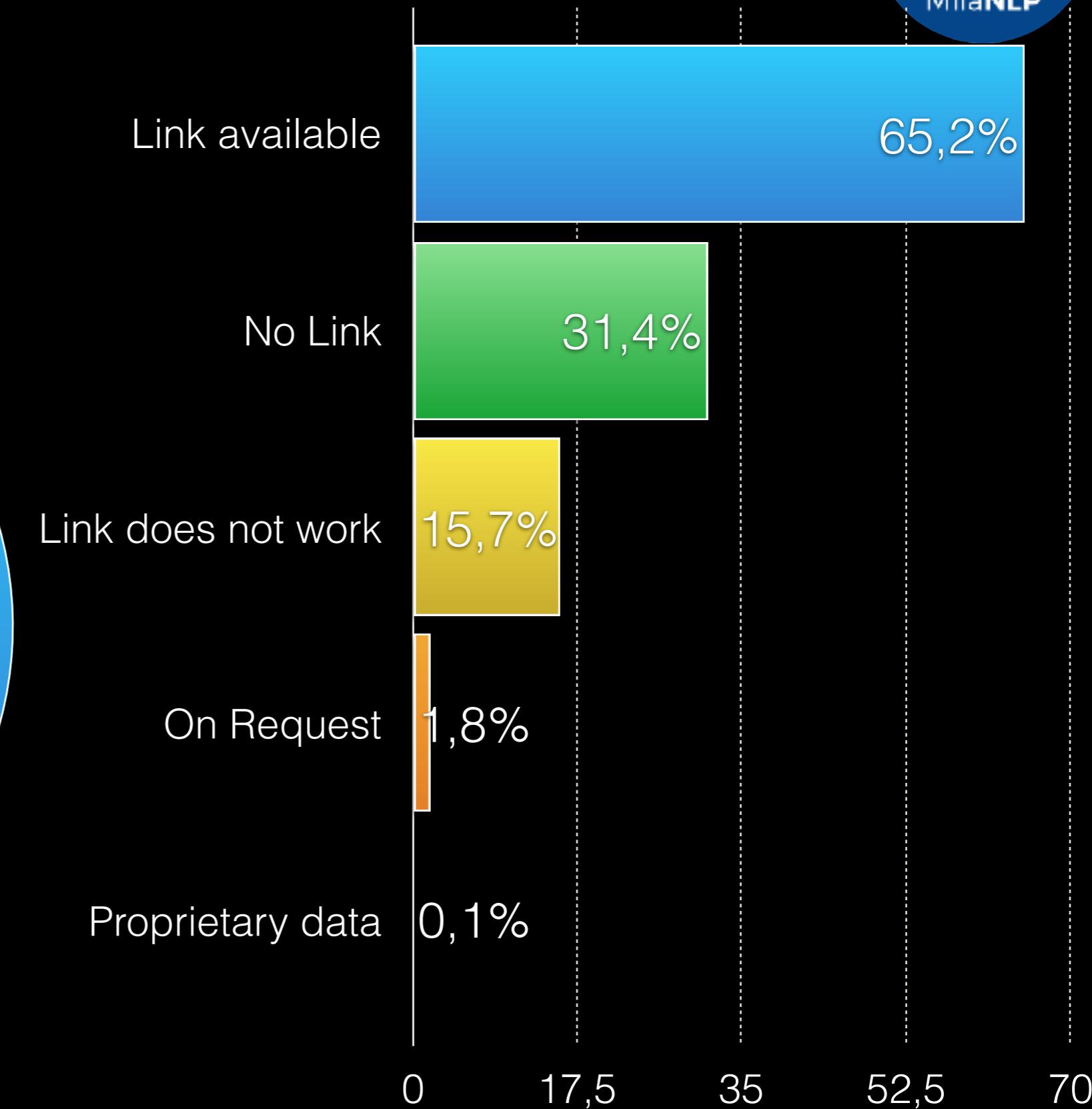
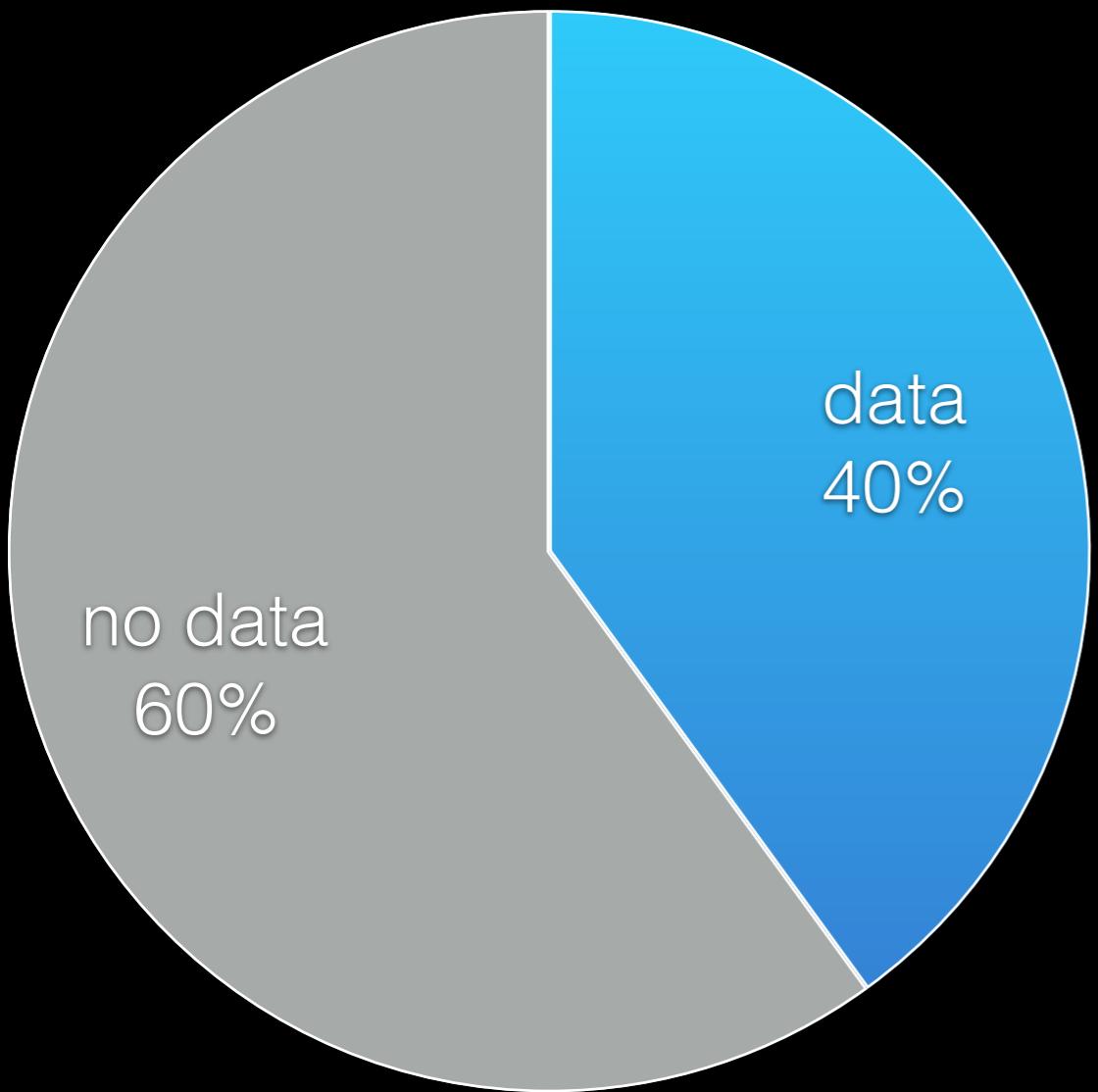


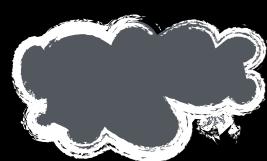
MAKE DATA AND
FEATURES EXPLICIT!





Replicability: Data





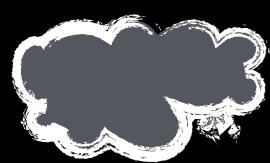
The Bender Rule



"Do state the name of the language that is being studied, even if it's English."

- Emily Bender





Idea!



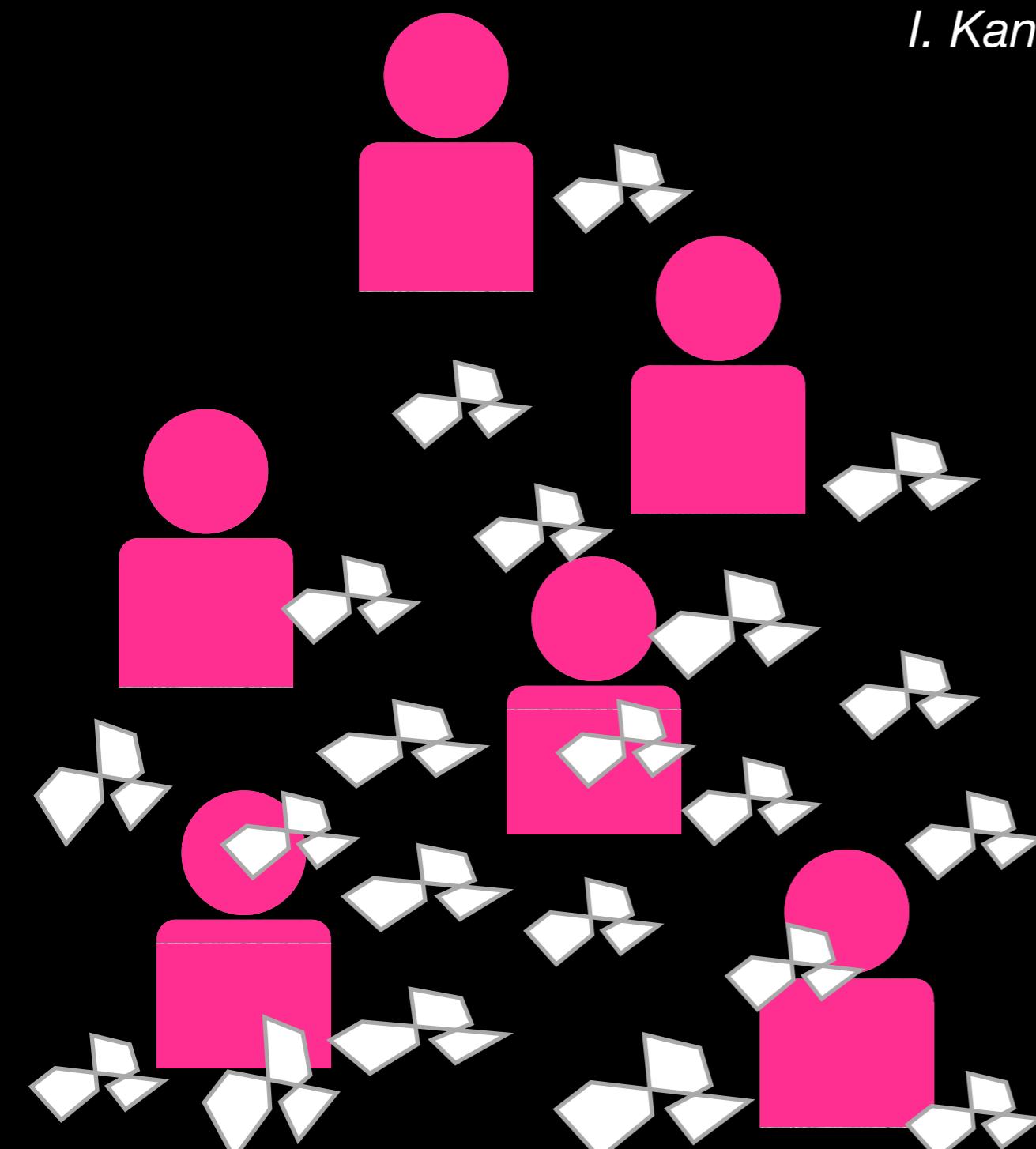
THINK ABOUT IMPLICATIONS



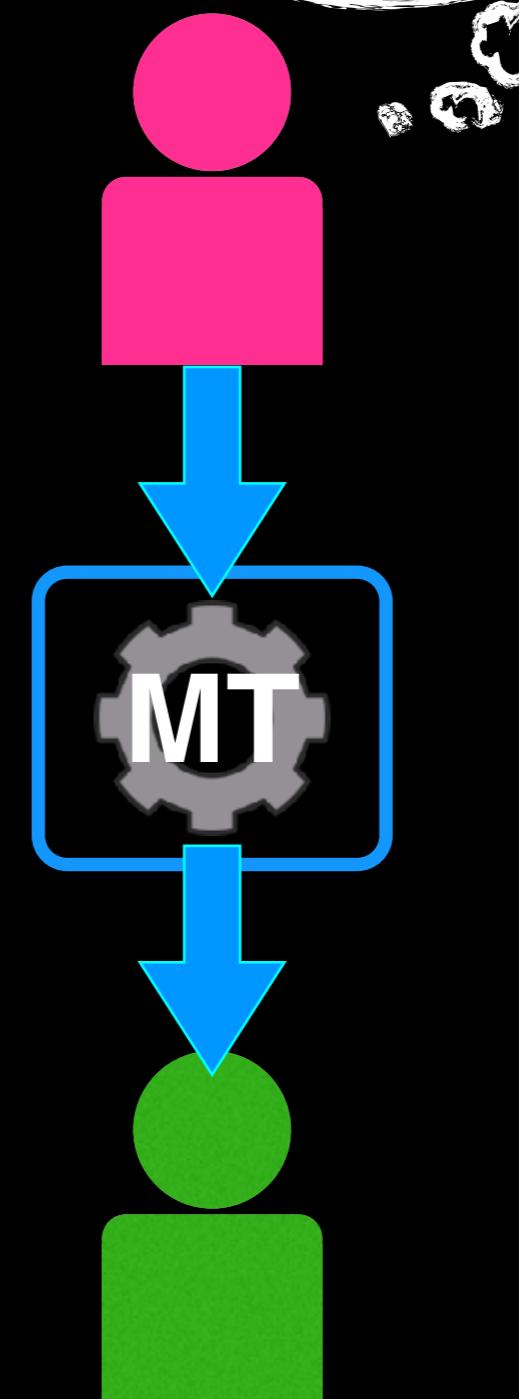
Generalization and Autonomy

"Act only according to that maxim whereby you can, at the same time, will that it should become a universal law"

I. Kant



THAT'S NOT
WHAT I SOUND LIKE!



Technology as Social Experiment

How can I make it safe?

Do my subjects know they're
in it?

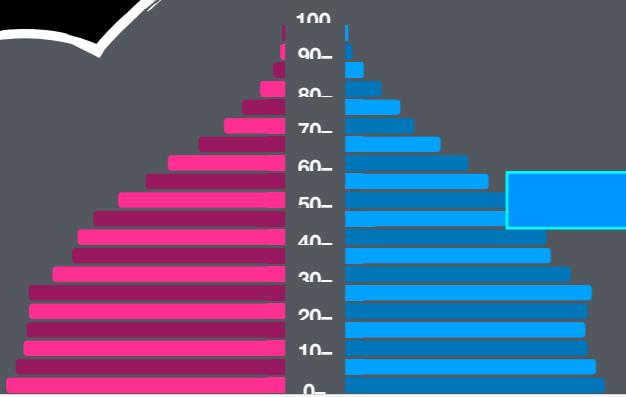
Can they control it, opt out?



Wrapping Up

Bocconi

Sources of Bias



SELECTION



ANNOTATION

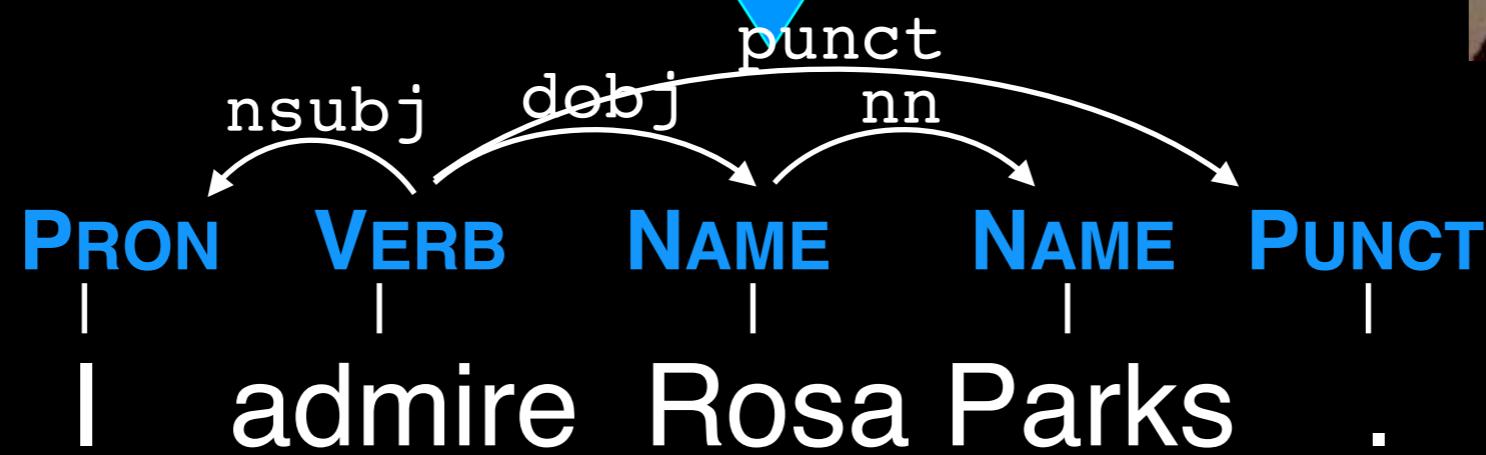


EMBEDDINGS

MODELS



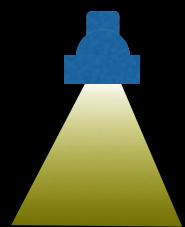
DESIGN



Bocconi

What can we do?

Source	Problem	Countermeasures
	Exclusion	better collection, post-stratification, priors
	Label Bias	better training, annotation models, disagreement weighting
	Overgeneralization	dummy labels, error weighting, adversarial learning
	Exposure	document, consider possible impact, educate



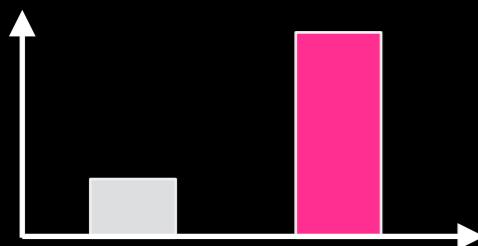
Tackling Bias



Fairness



Personalization



Performance

Take-home points



- Beware of **bias** from **data, annotations, embeddings, models, and design**
- Apply **countermeasures** where possible
- Know your models **will** be used in unintended ways
- Ask yourself:
"Am I comfortable with my system classifying me?"

www.dirkhovy.com/portfolio/papers



Thank you!



@dirk_hovy

www.dirkhovy.com

Bocconi