

Natural Language Processing

Lecture 15

Dirk Hovy

dirk.hovy@unibocconi.it

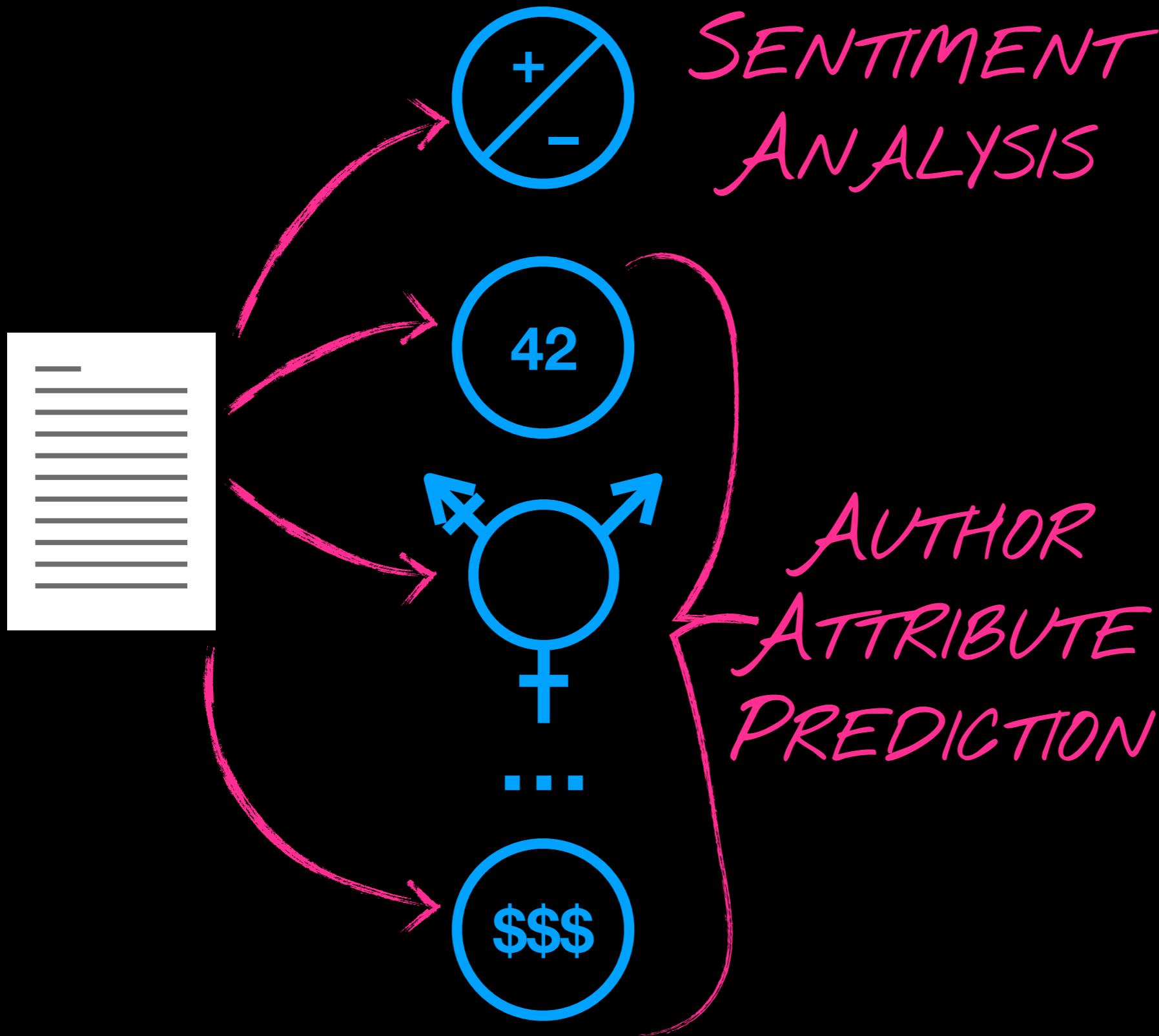
 @dirk_hovy

Bocconi

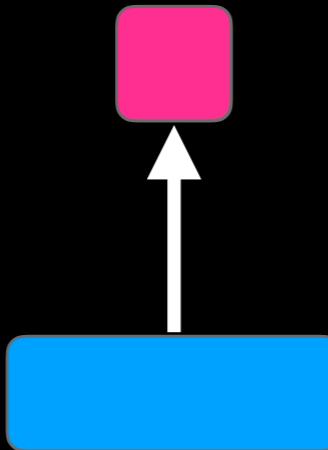
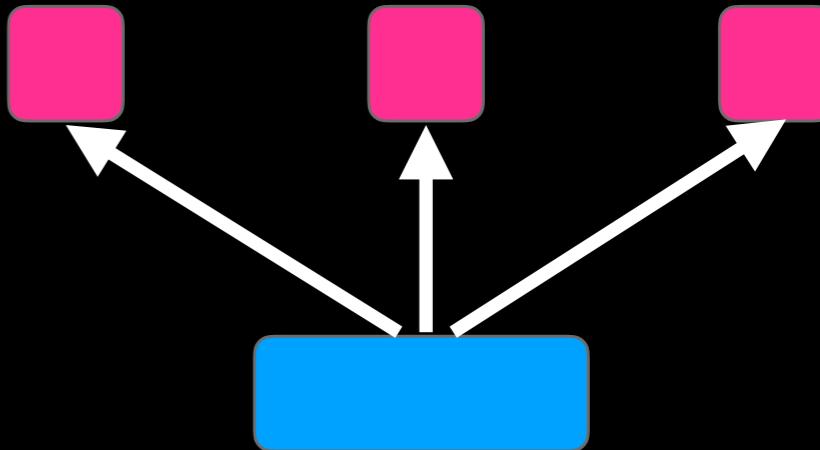
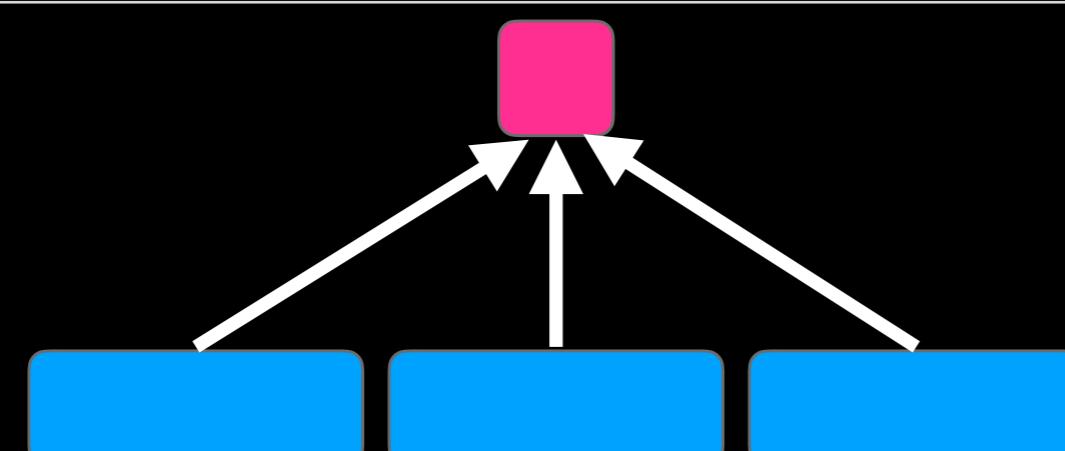
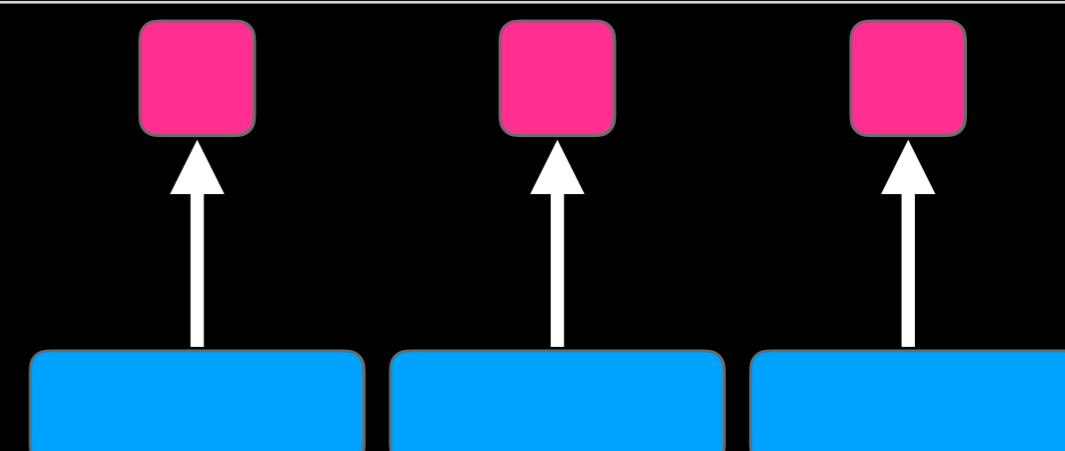
Text Classification



Examples



Types of Text Classification

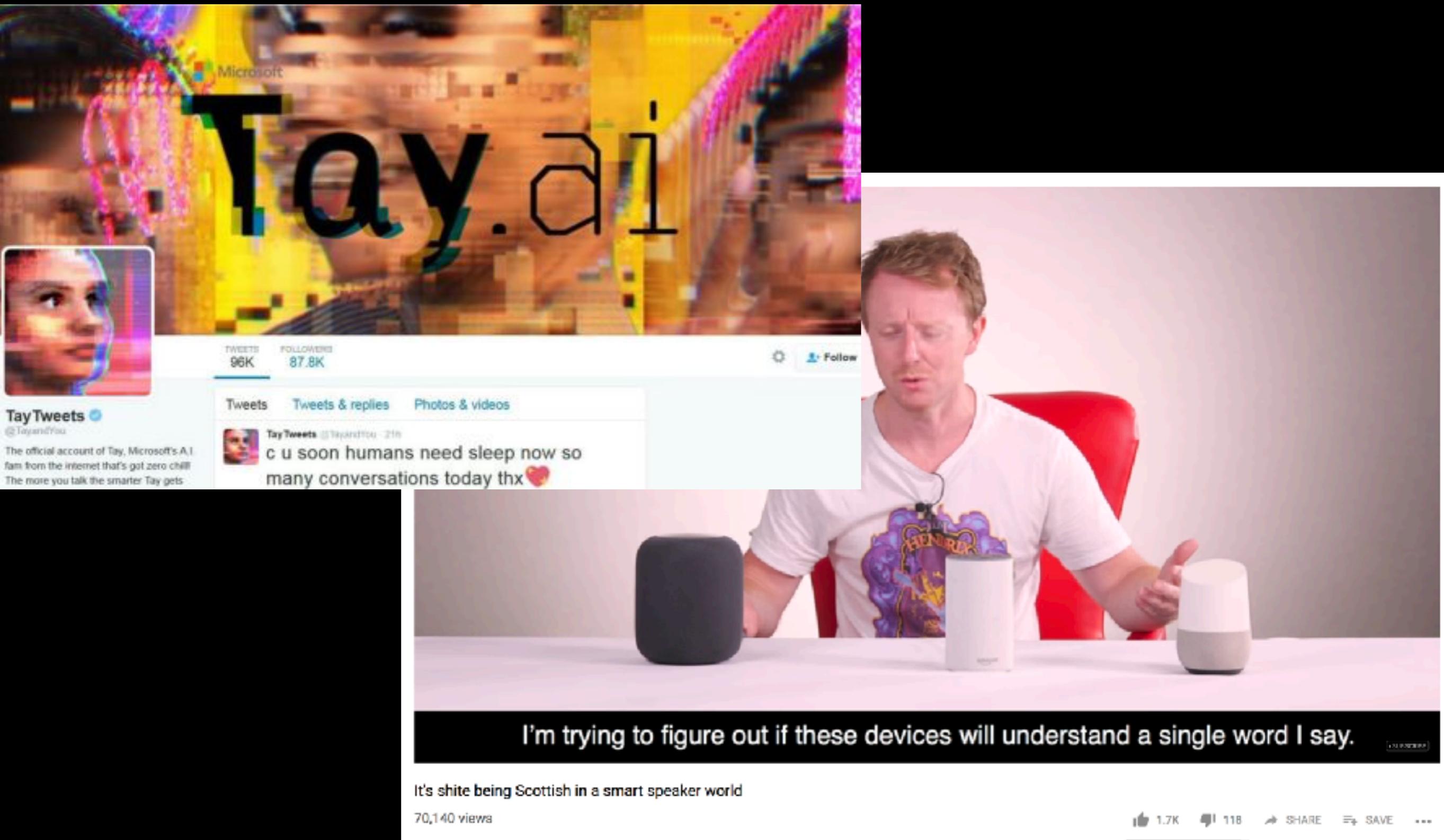
	Fixed length	Variable length
Fixed length	 <p>Logistic Regression, Perceptron, Feed-Forward Network, Random Forest, Naive Bayes, SVM, ...</p>	 <p>Multitask Learning, Decoder</p>
Variable length	 <p>Convolutional Neural Networks (CNN)</p>	 <p>Recurrent Neural Networks (RNN), Hidden Markov Models (HMM), Conditional Random Fields</p>

Goals for Today

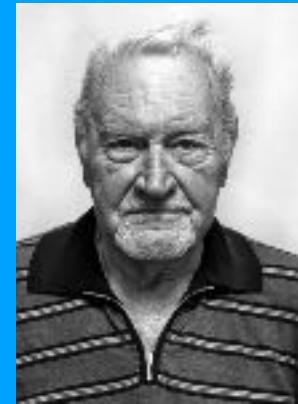
- Learn about **ethical issues**
- Understand potential **sources of bias**
- Learn about **counter measures**

Ethical Considerations

Biased Language Systems



Language Biases

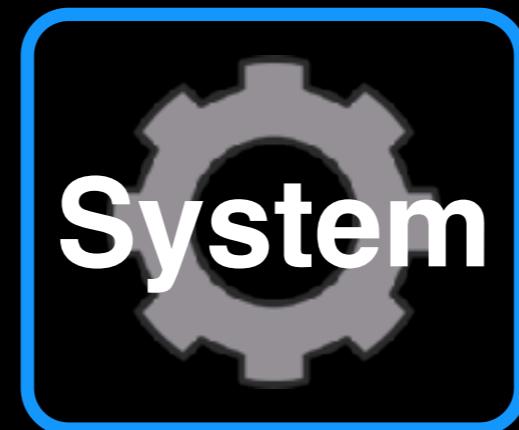


Example 1

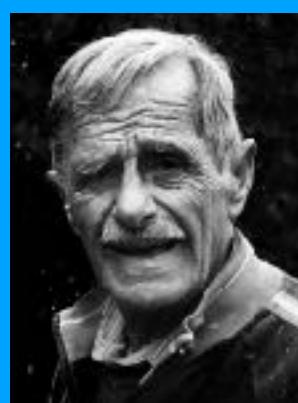
I don't understand you...



Example 2



Hello,
computer

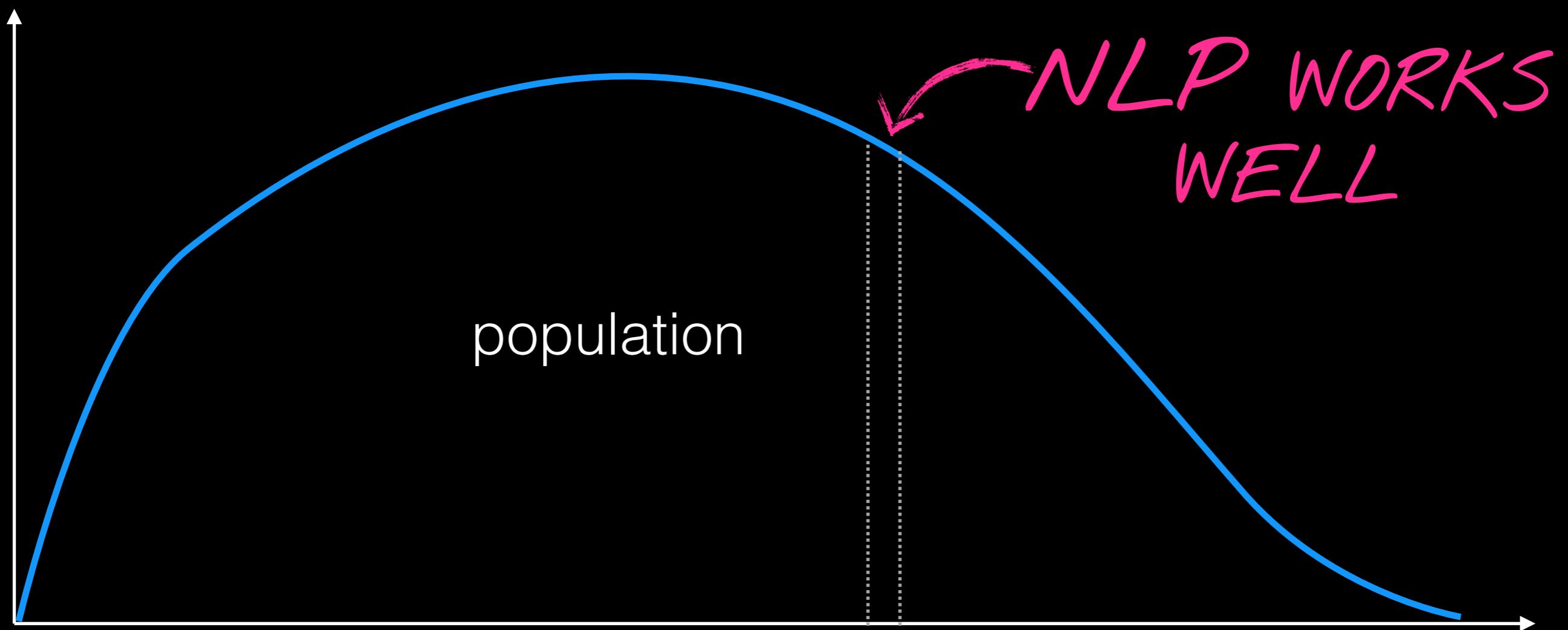


Example N

Shite...



The Consequences



Solutions?

SYMPTOMS



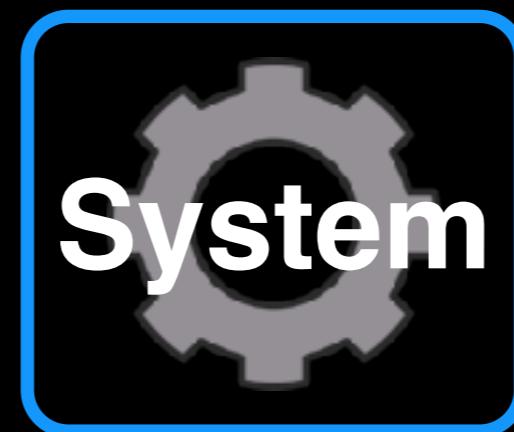
Example 1



Example 2



Example N



CAUSE



Language

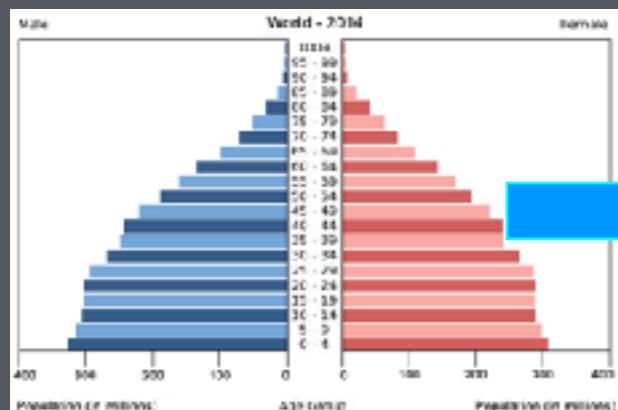
Who cares about bias?

Philosophers

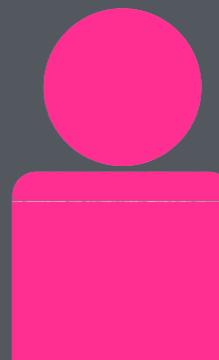
WE ALL SHOULD...

Legal

Sources of Bias



SELECTION

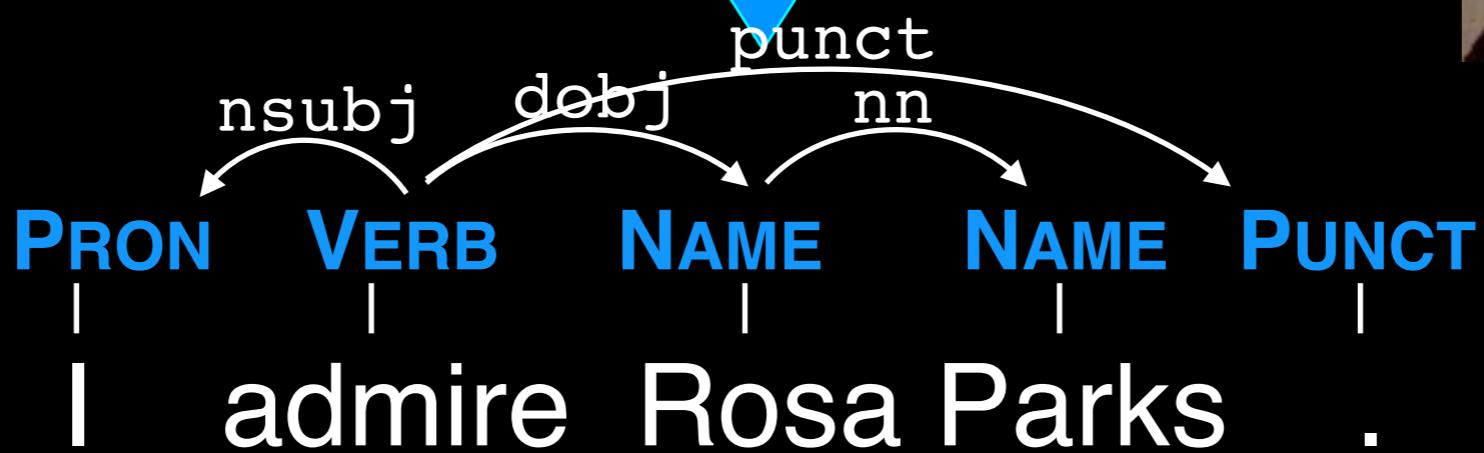


ANNOTATION

MODELS

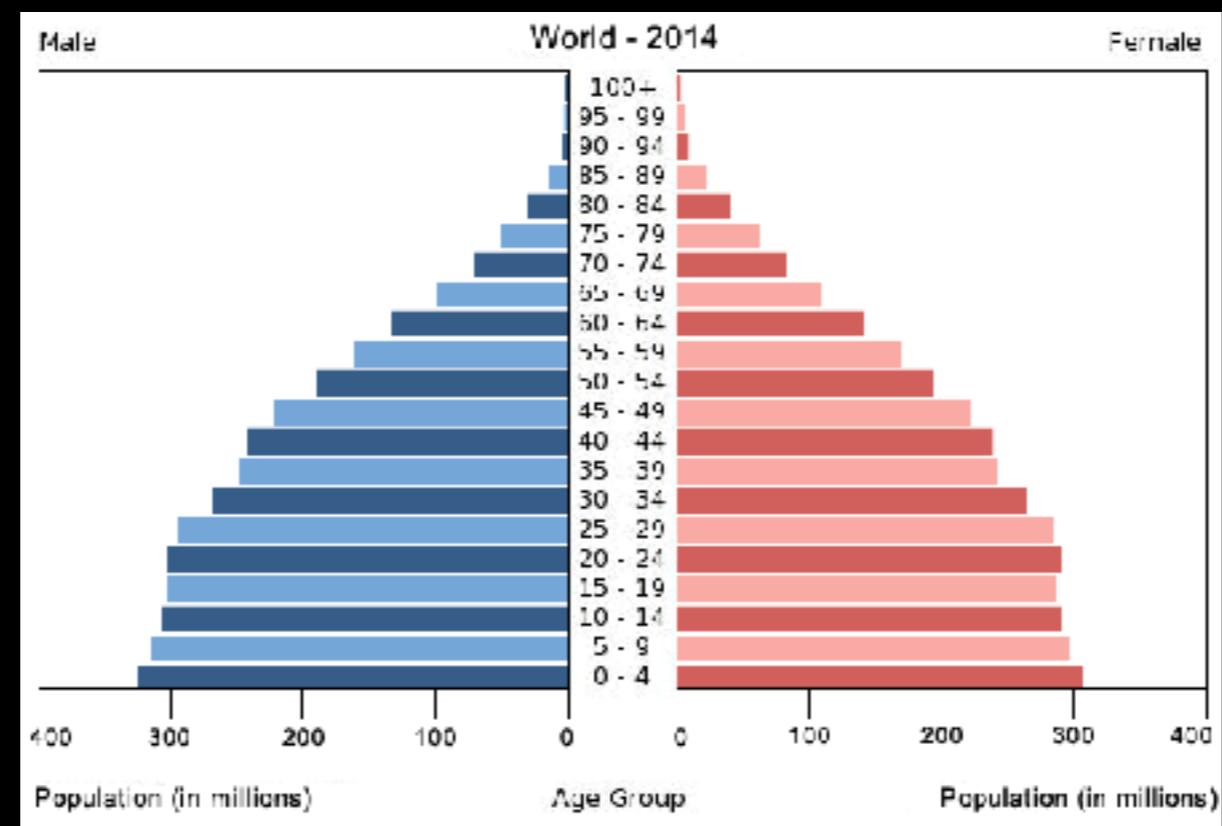


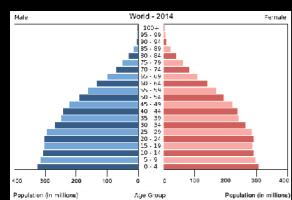
DESIGN



Bocconi

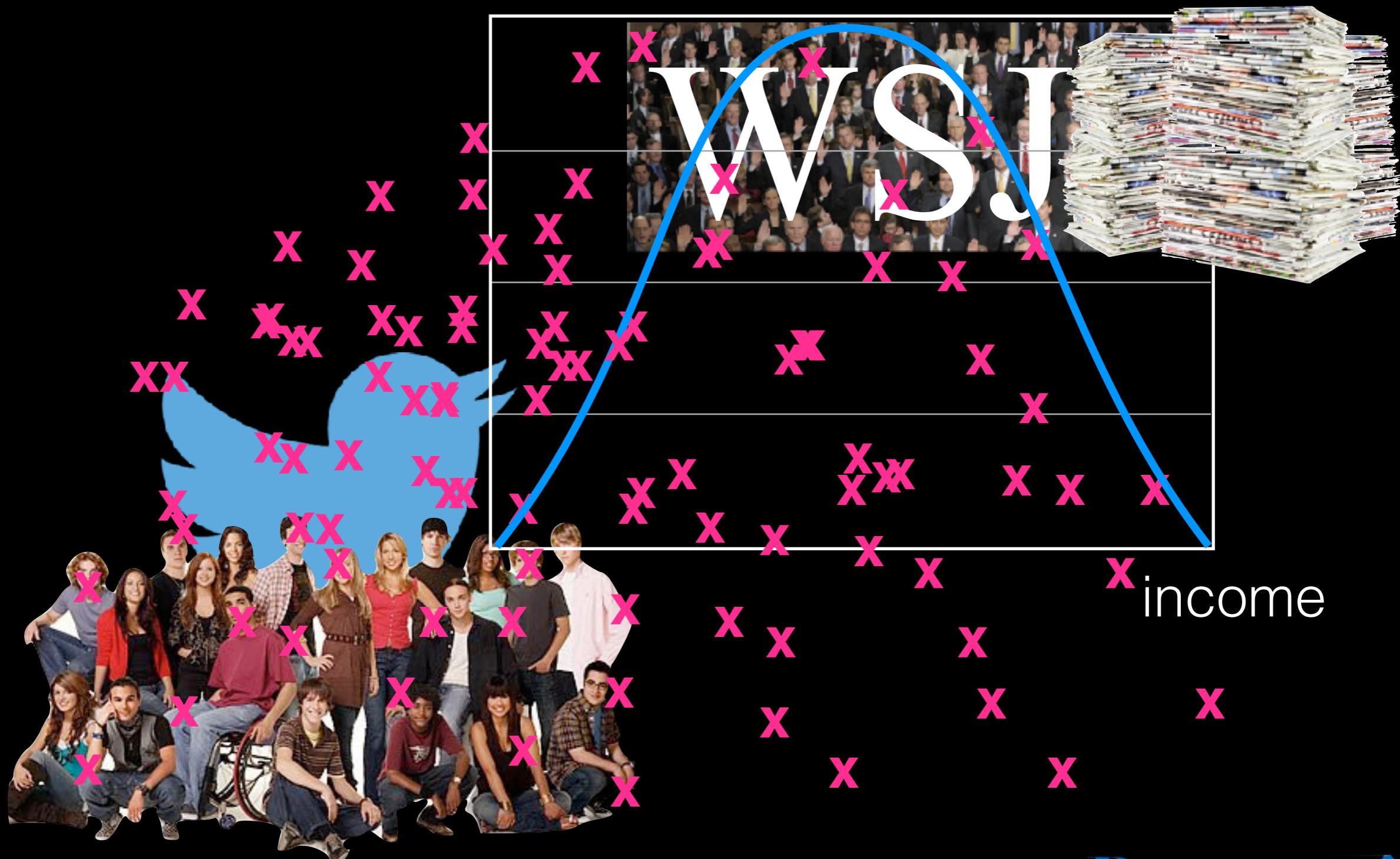
Data Bias

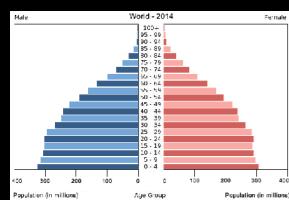




Distributions

age

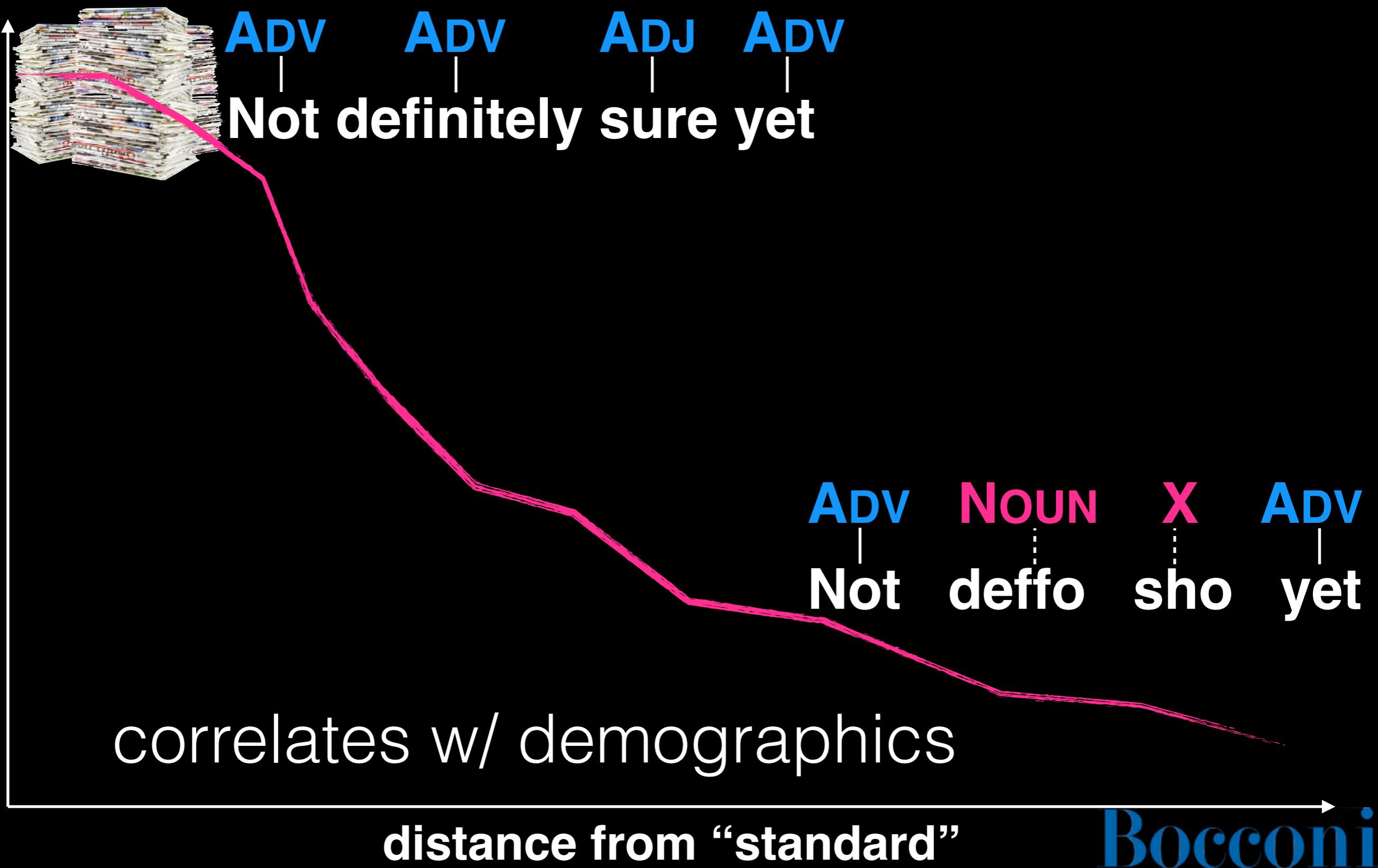


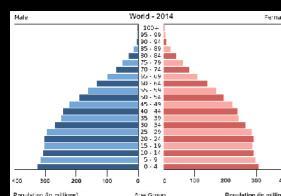


Jørgensen et al. (WNUT 2015)
Høvy & Søgaard (ACL 2015)

NLP
performance

The WSJ Effect

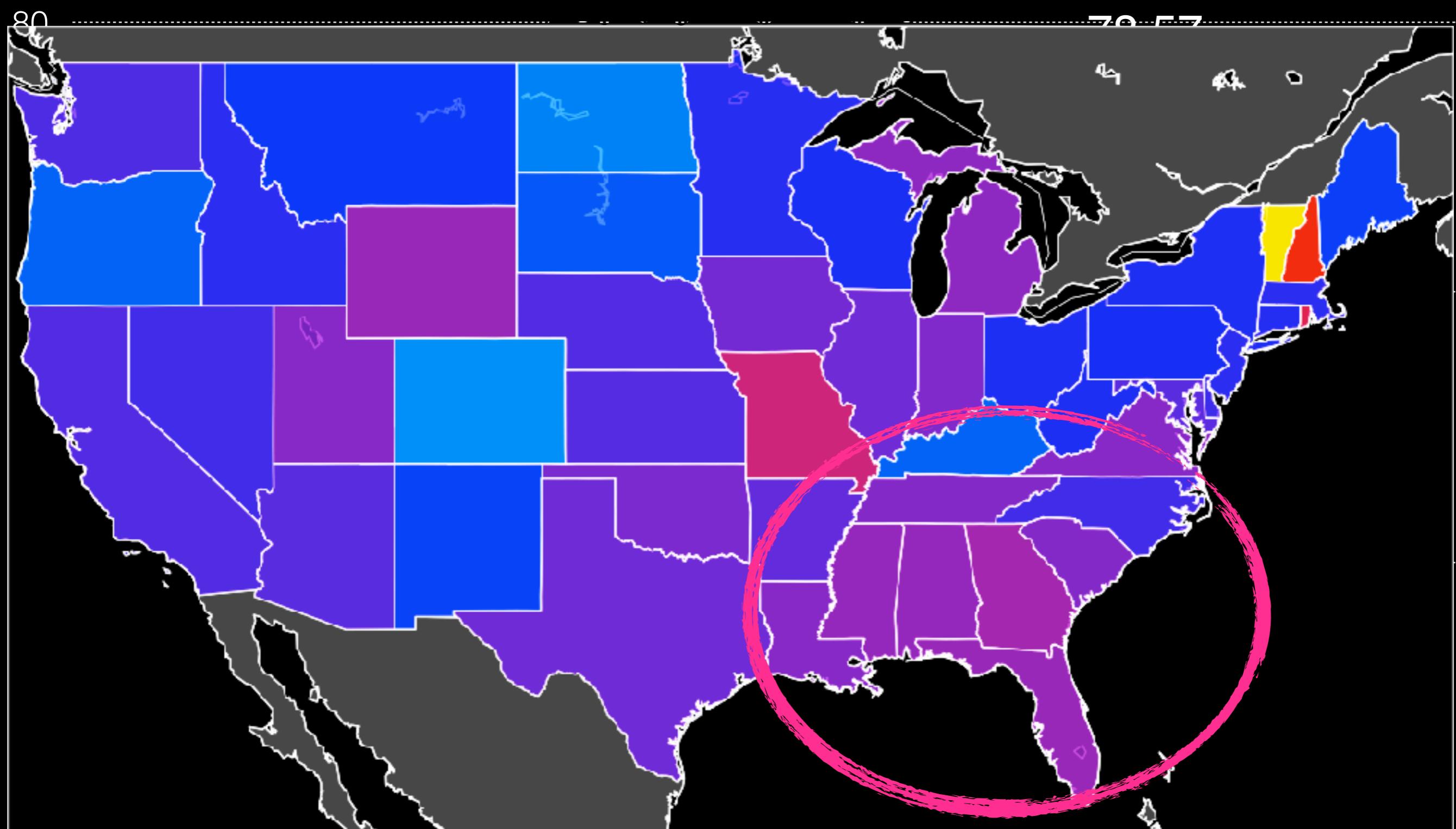


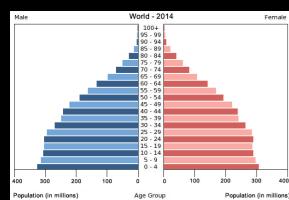


Jørgensen et al. (WNUT 2015)
Hovy & Spruit (ACL 2016)

Exclusion

F1





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

Exclusion

accuracy

100

Male

World - 2014

Female

ing
views
ned

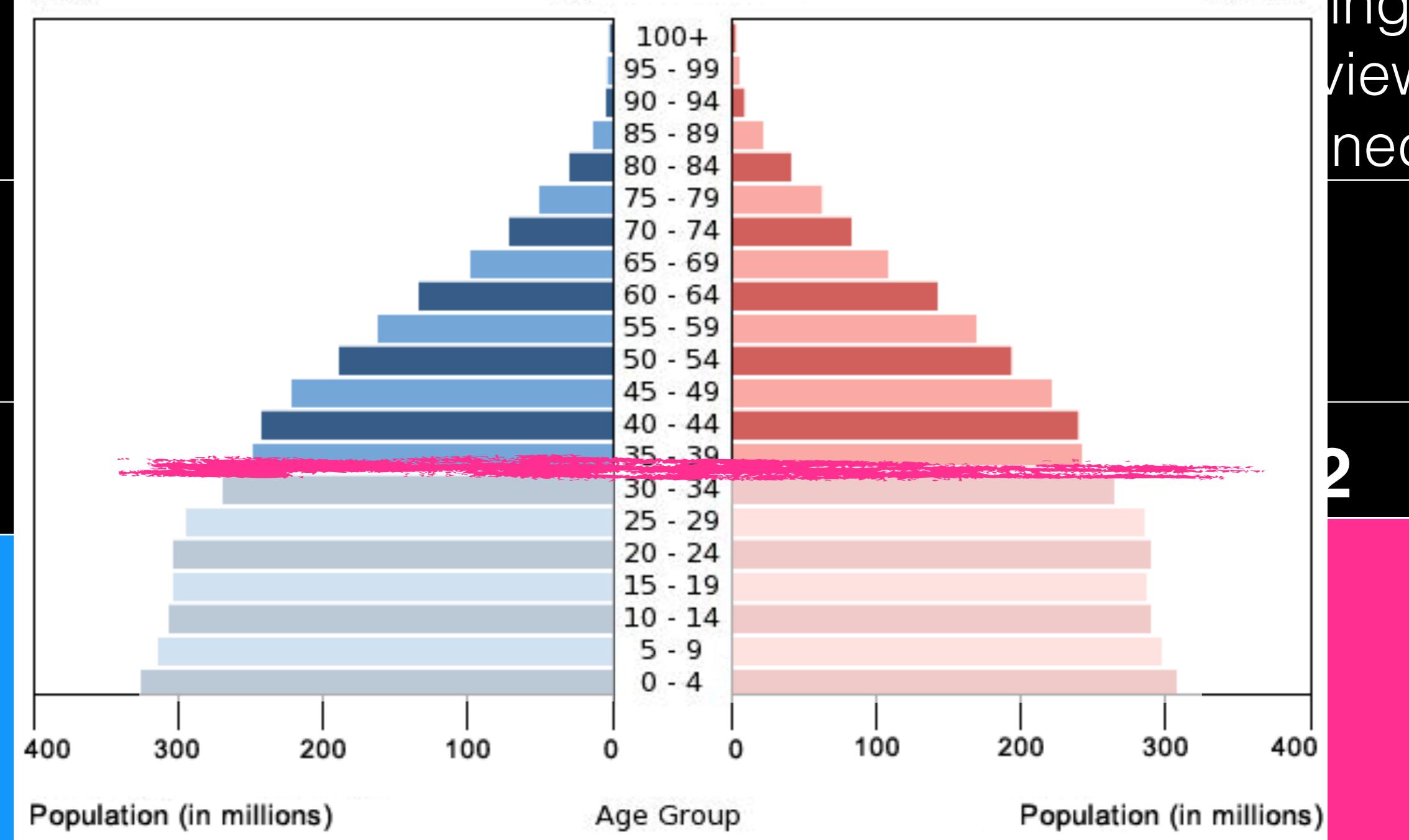
95

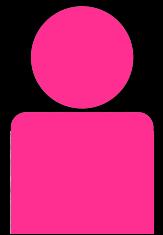
90

85

80

2





Annotator Bias



It's a
particle!



No! It's an
adposition!

PRON VERB

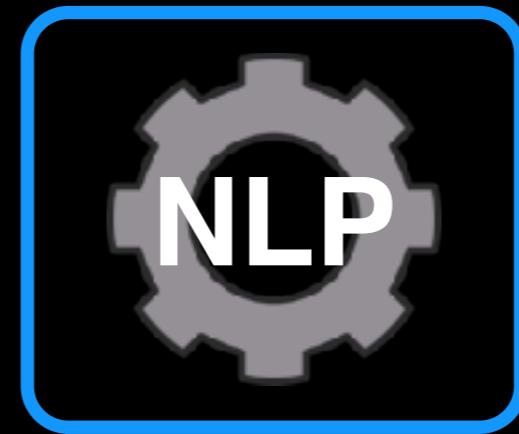
PRON VERB

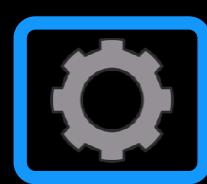
PRT NOUN NUM

ADP NOUN NUM

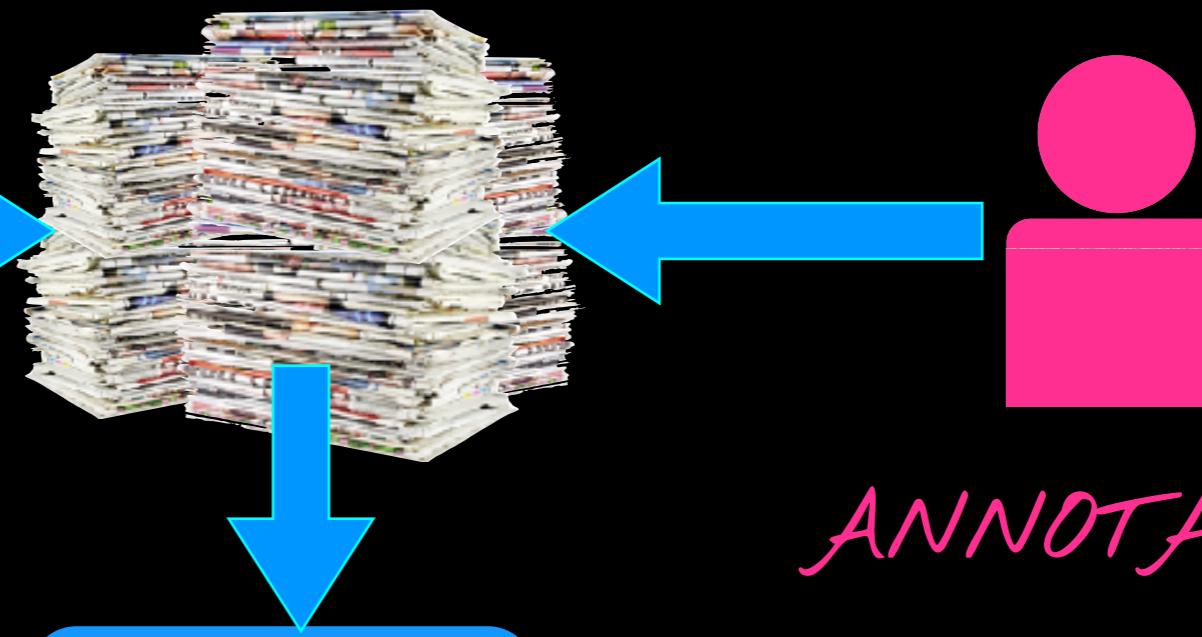
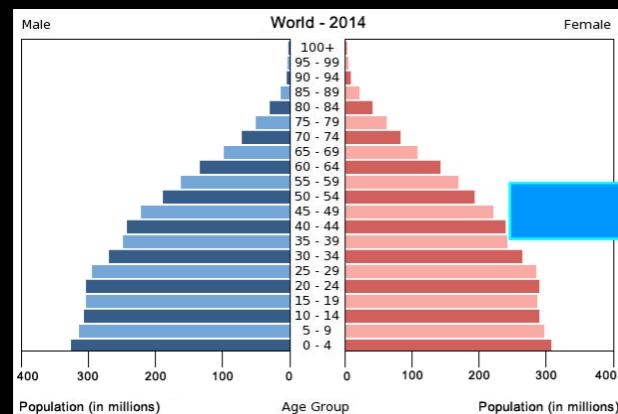
it comes out apr 30

Model Bias





Biased Models

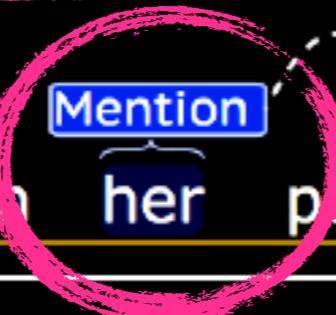


Wrong Coreference

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 !



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**

Models Amplifying Bias

BIAS = 0.66



Agent: WOMAN



Agent: MAN



Agent: WOMAN

BIAS = 0.84



Agent: WOMAN



Agent: WOMAN



Agent: WOMAN



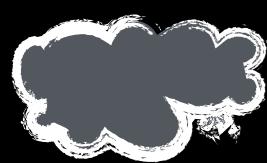
Agent: MAN



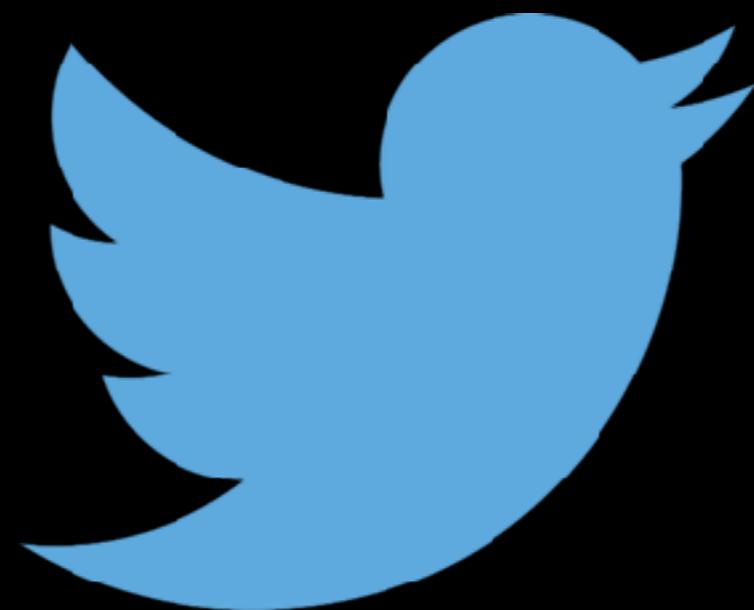
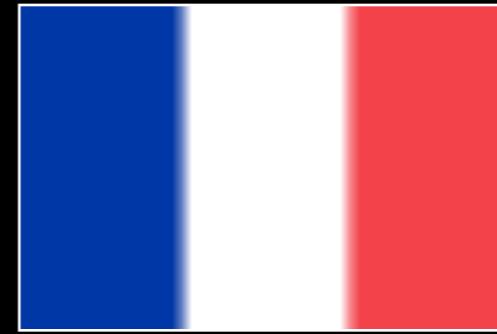
Agent: WOMAN

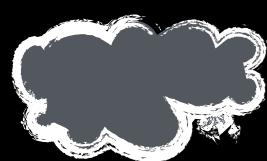
Design Bias





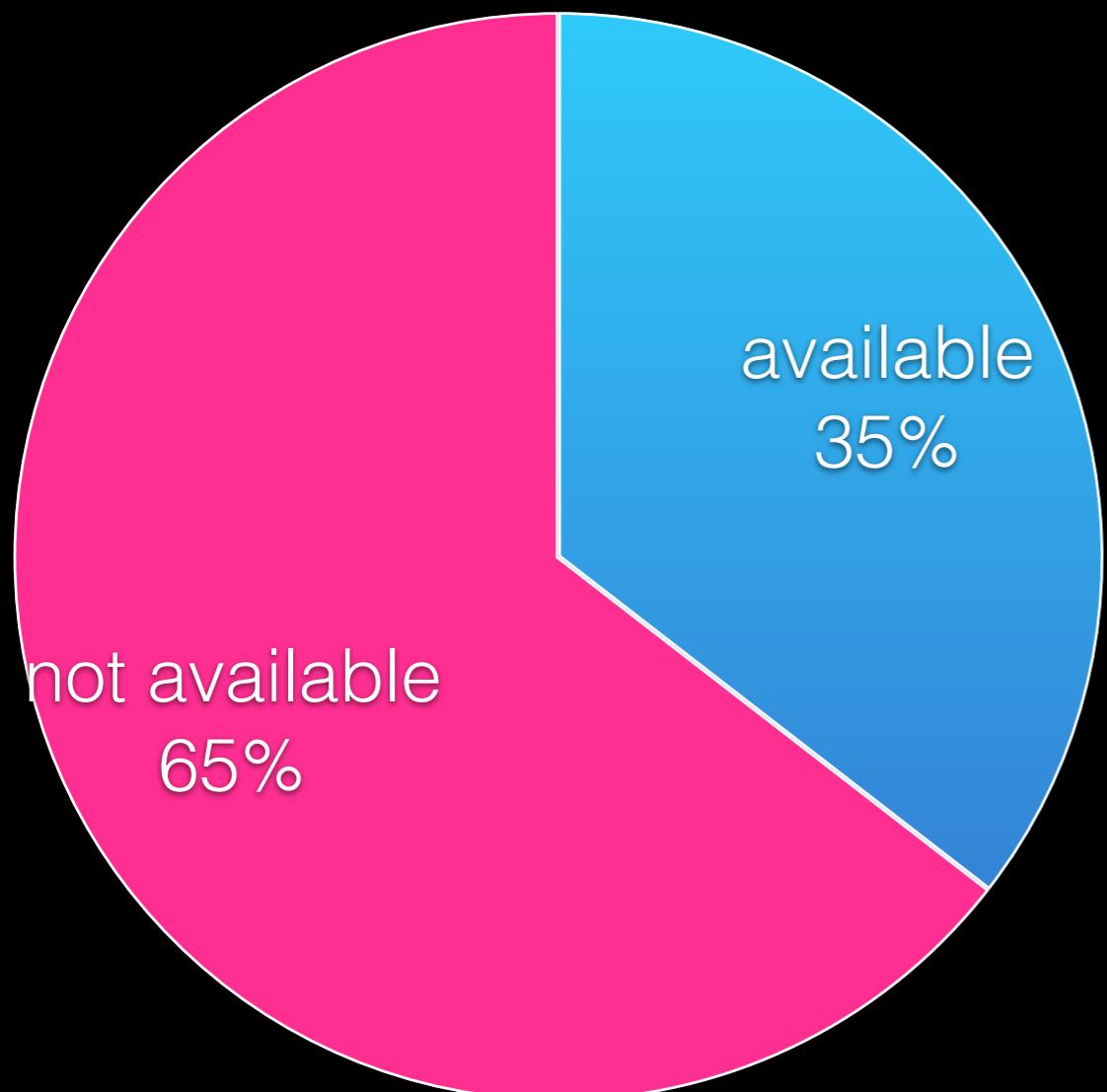
Exposure



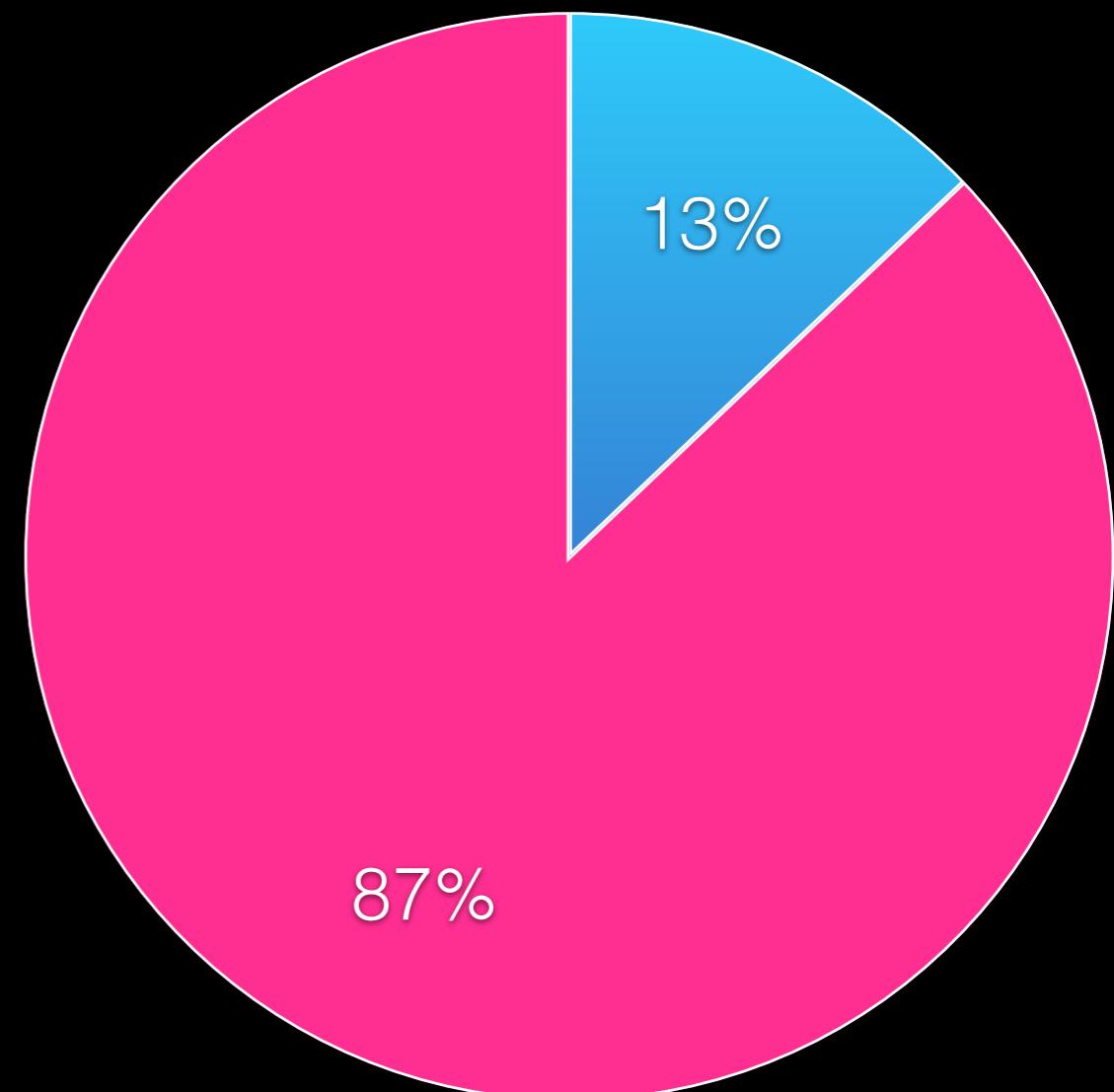


Under-Exposure

treebanks



semantic resources



evaluation

Over-Exposure



American
New York City
English

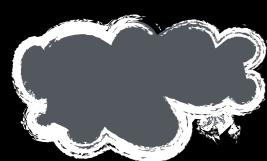


Nigagios
English
16m

POS tagging

Discourse

Bocconi



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

English Turkish Spanish Detect language ▼

She is a doctor.
He is a nurse.

31/5000

English Turkish Spanish ▼

O bir doktor.
O bir hemşire.

Translate

English Turkish Spanish Turkish - detected ▼

O bir doktor.
O bir hemşire

NORMATIVELY WRONG

27/5000

English Turkish Spanish ▼

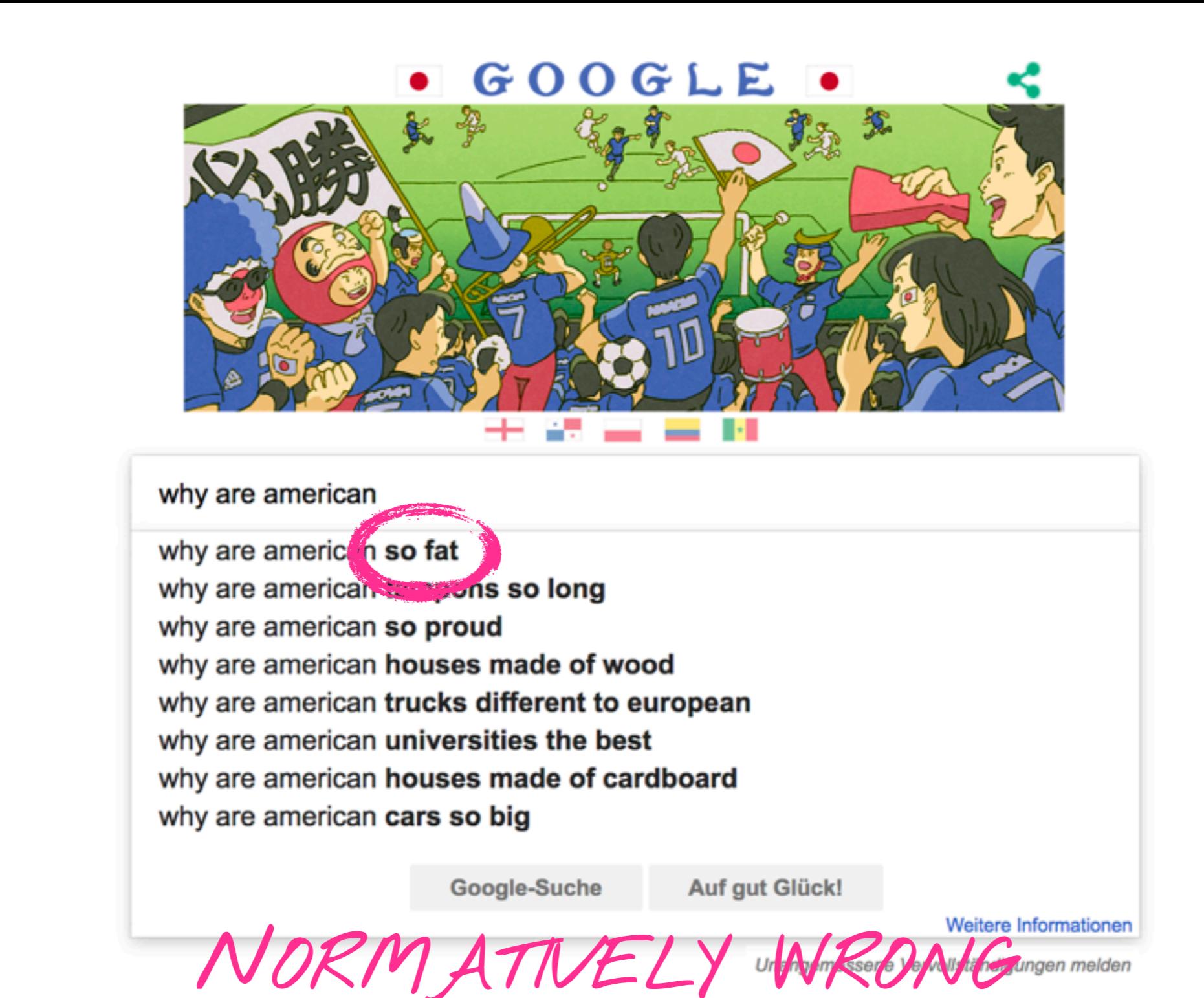
He is a doctor.
She is a nurse ✓

DESCRIPTIVELY WRONG

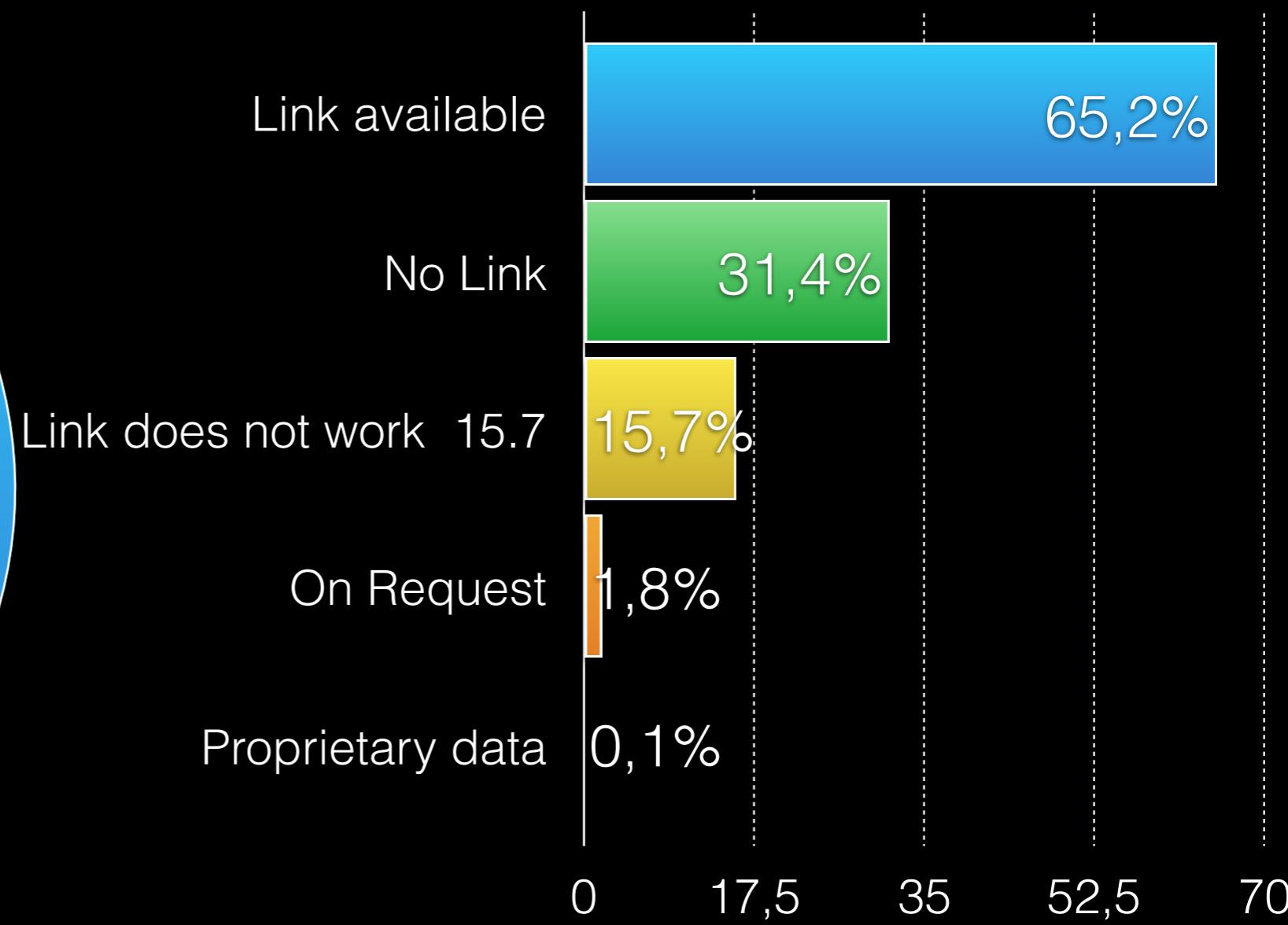
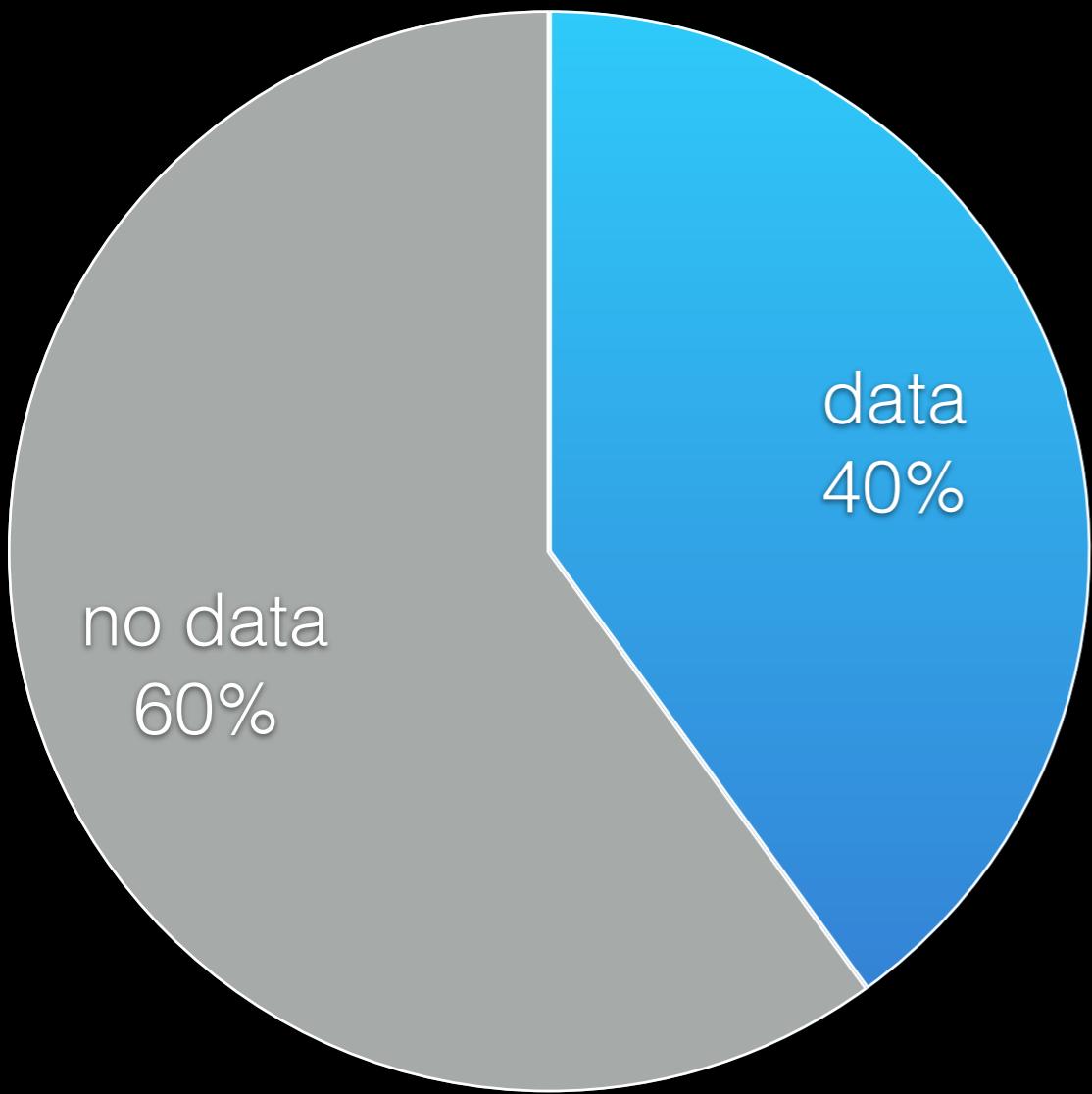
27/5000

Translate

Normative vs Descriptive Ethics

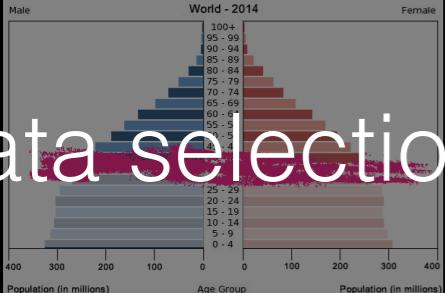


Replicability



Counter Measures

What can we do?

Source	Problem	Countermeasures
	Exclusion	stratification, priors
	Label Bias	annotation models, disagreement weighting
	Overgeneralization	dummy labels, error weighting, adversarial learning
	Exposure	always consider possible impact

Reducing Bias

$$BIAS = 0.66$$



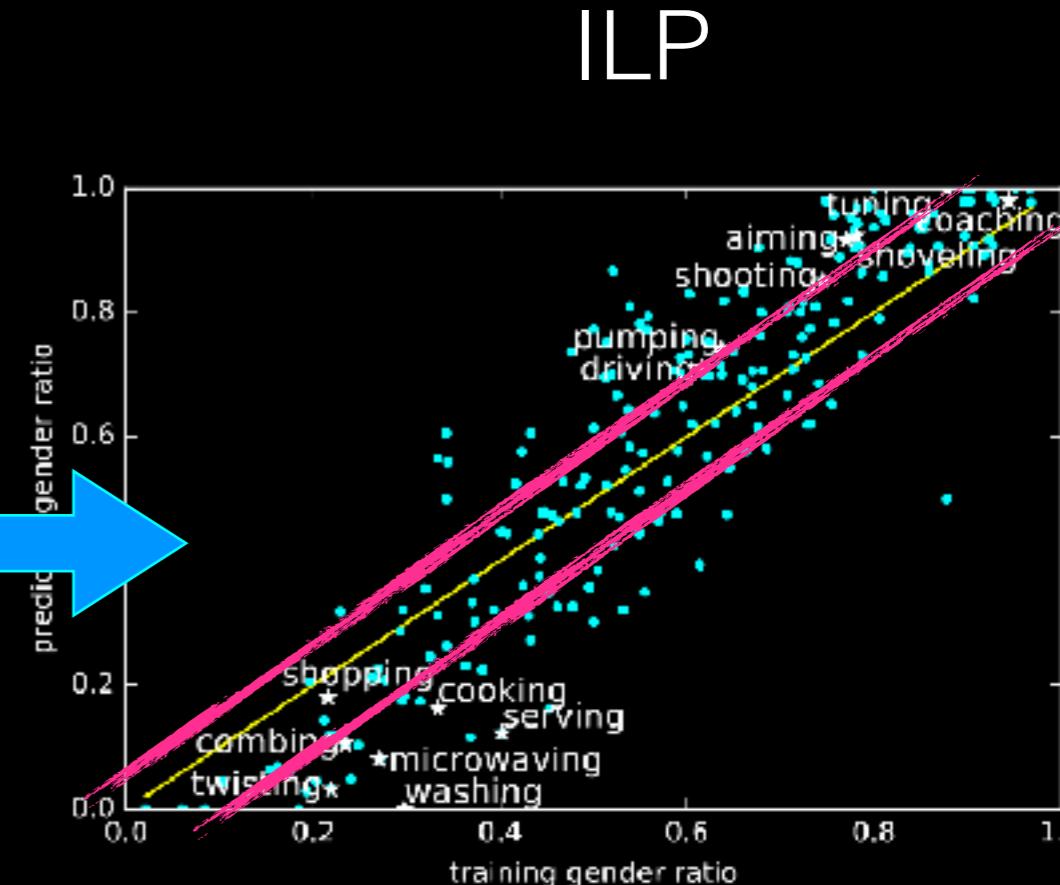
Agent: WOMAN

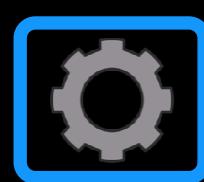


Agent: MAN

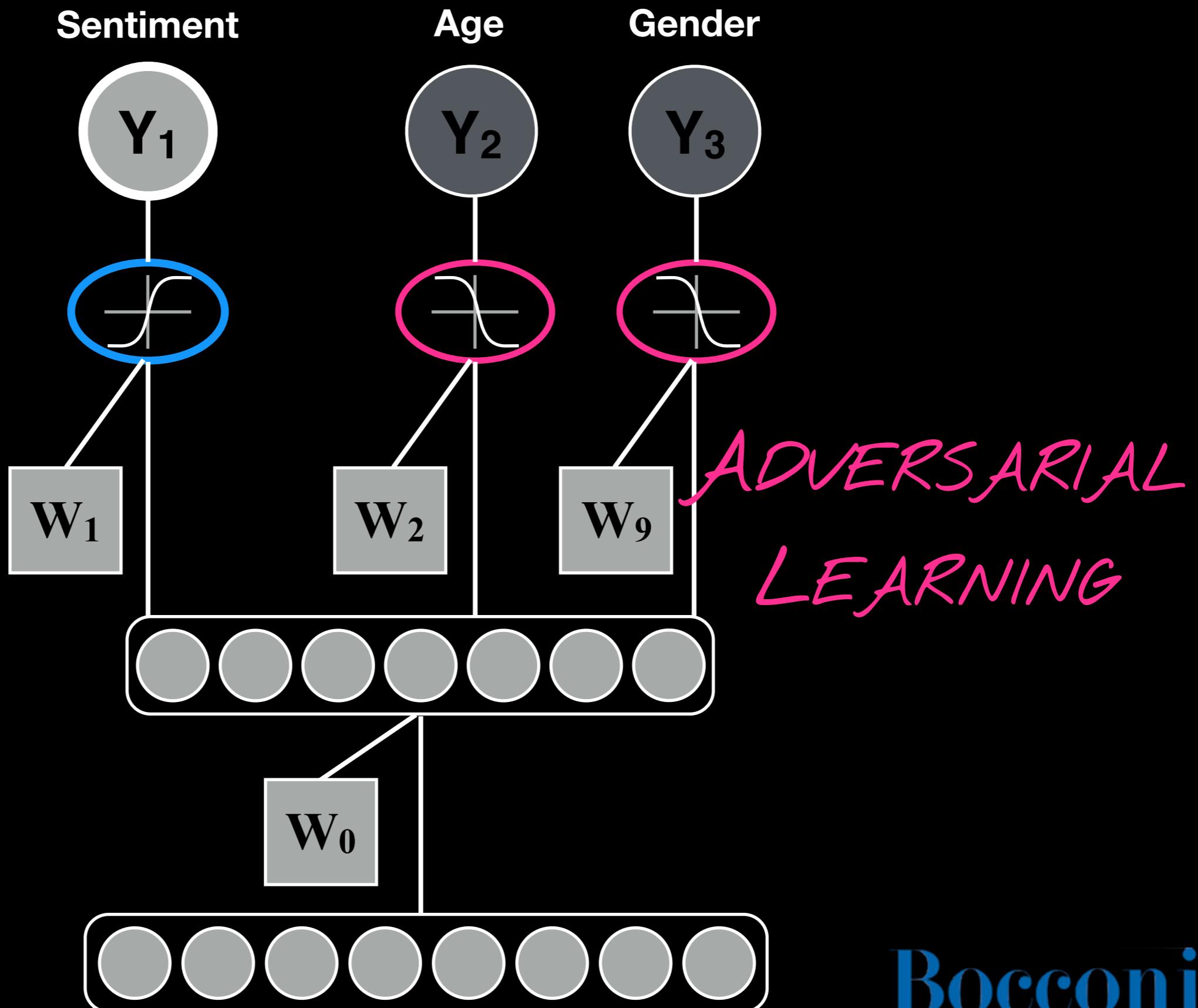


Agent: WOMAN

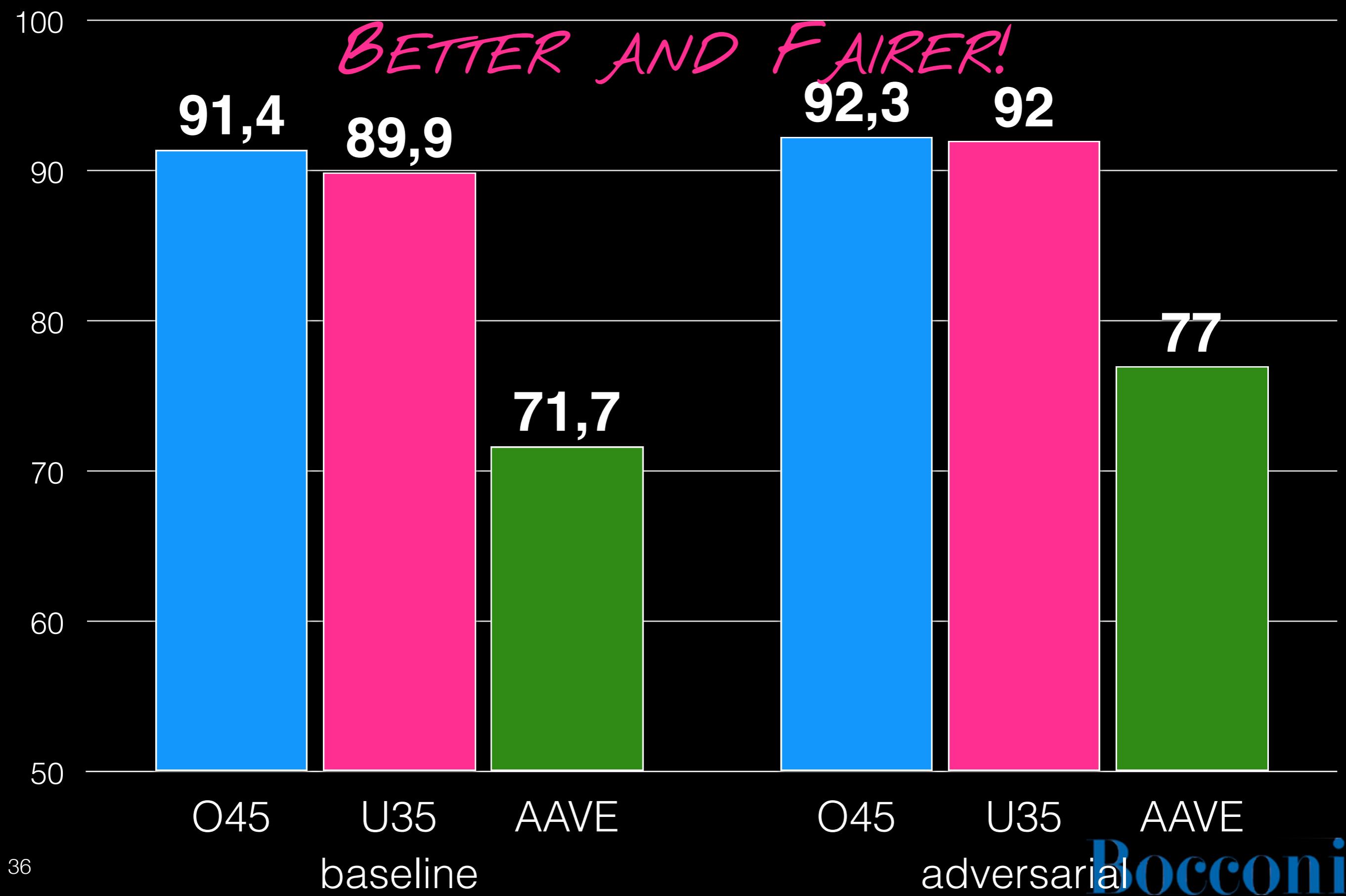




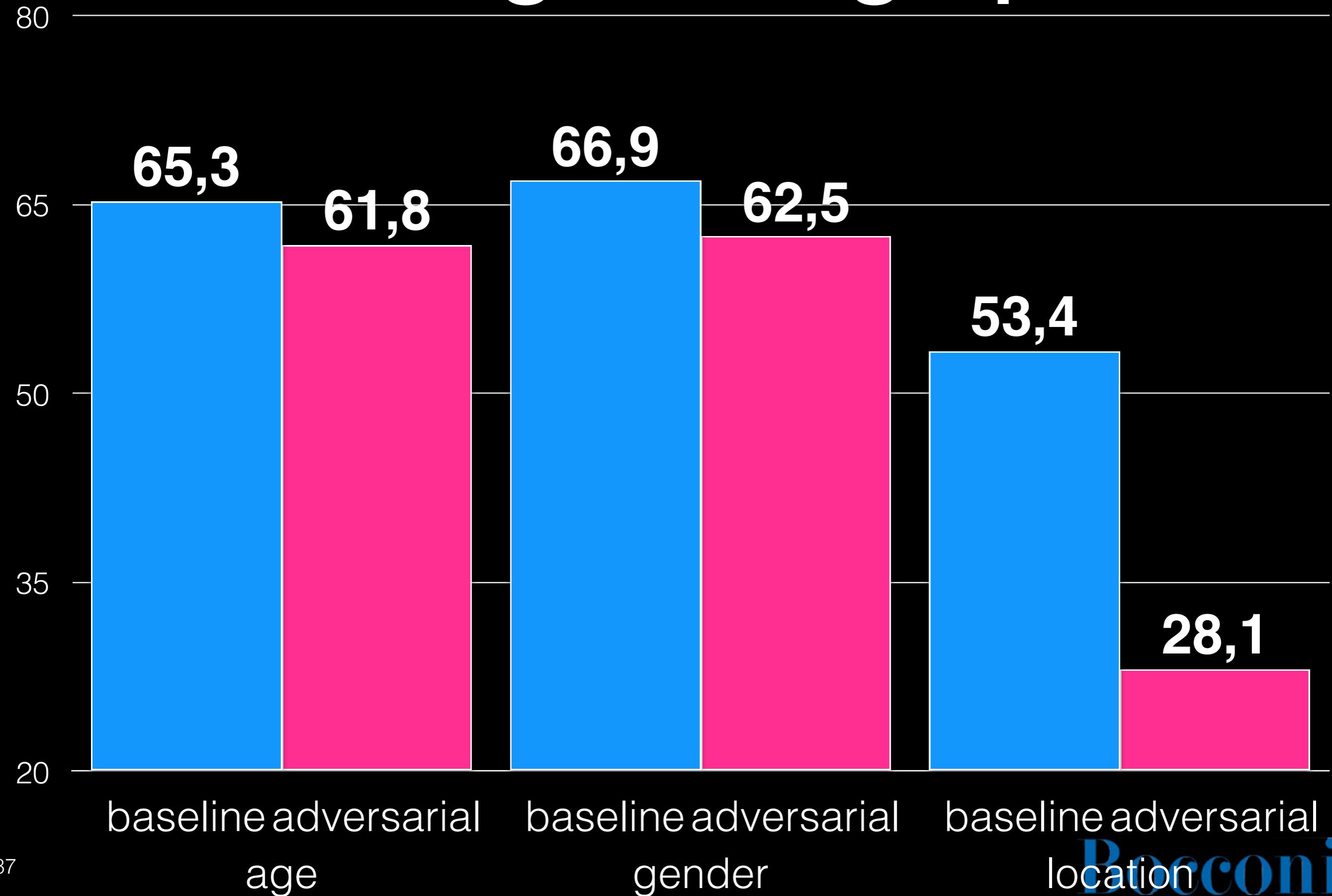
Adversarial Model



Results

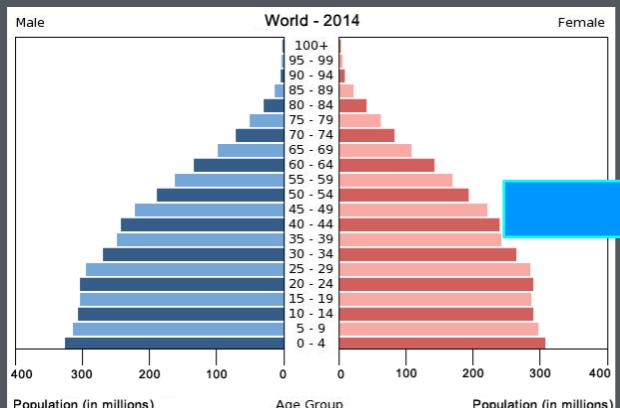


Protecting Demographics



Wrapping Up

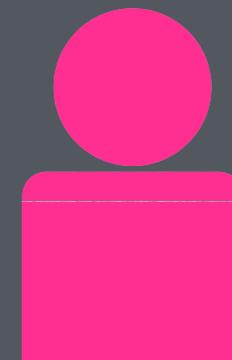
Sources of Bias



SELECTION



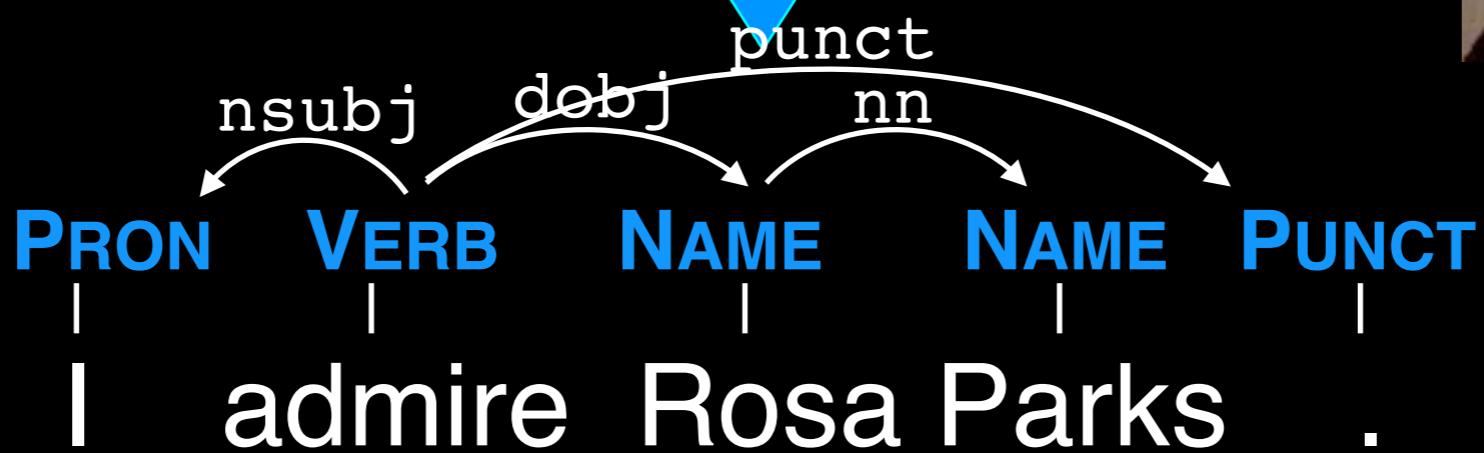
ANNOTATION



MODELS



DESIGN



Bocconi

Open Questions

- Carrots or sticks (GDPR)?
- Do we have to be more humble? How? Where?
- What about normative vs. descriptive ethics?
- Consequences for privacy, security, and user dignity?
- What if you have to choose between ethics and performance?
- Are we asking a lot from small companies?
- Research into bias creates biased data sets. Does it have to?
- Where does it stop: what about unknown sources of bias?

Take-home points

- Beware of **bias** from **data, models, or design**
- Apply **countermeasures** and check
- Ask yourself: "*Am I comfortable having my system classify myself?*"