



Five Sources of Biases in NLP

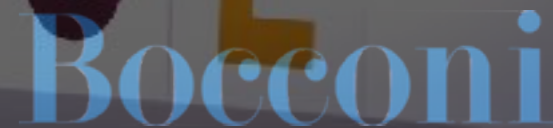
(and What to Do about Them)

IT'S NOT
JUST DATA

Dirk Hovy
Bocconi University, Milan

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

 @dirk_hovy



Text Data is Exploding



Exabytes = 1M TB

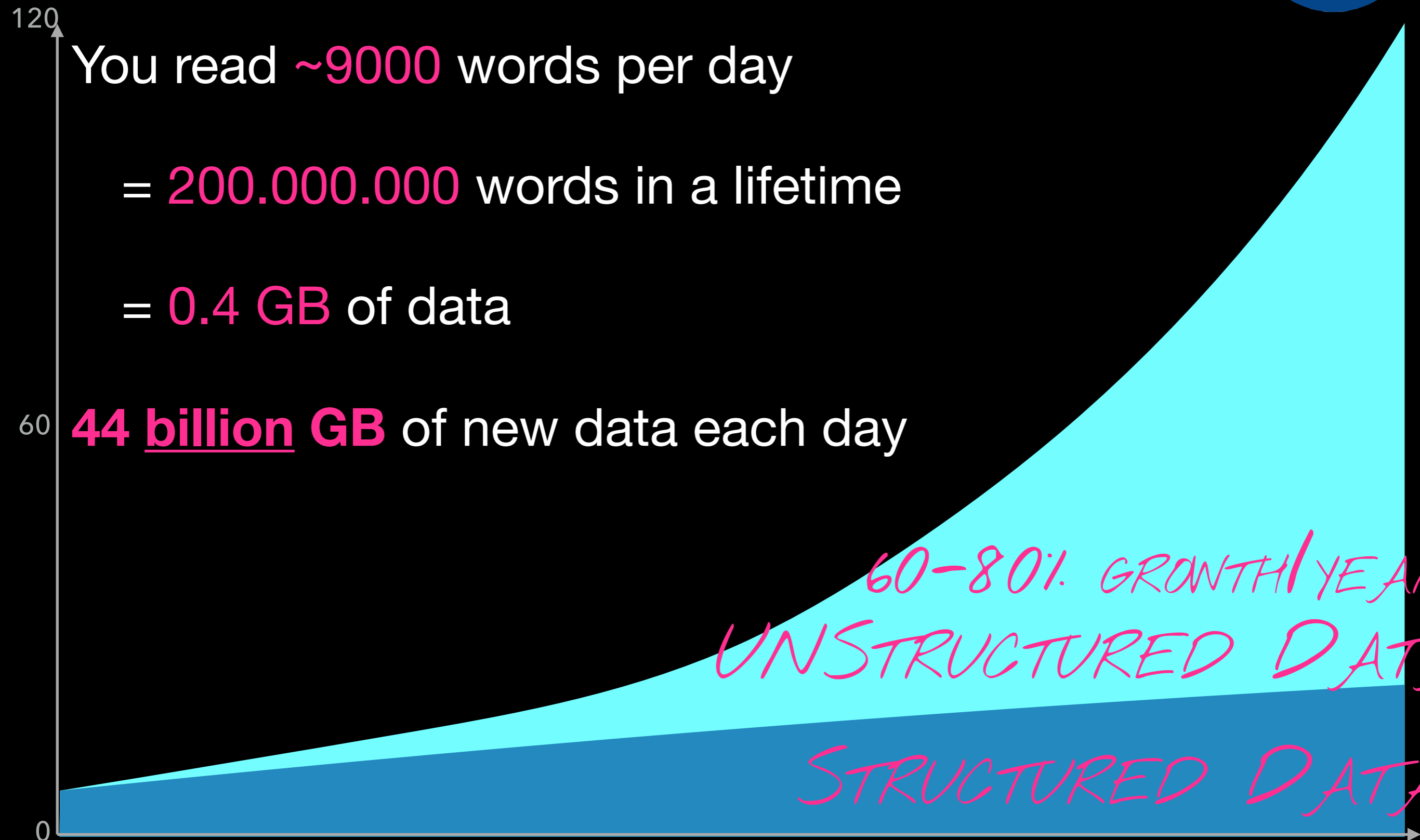
You read ~9000 words per day

= 200.000.000 words in a lifetime

= 0.4 GB of data

44 billion GB of new data each day

60-80% GROWTH/YEAR
UNSTRUCTURED DATA
STRUCTURED DATA



NLP is booming



³ Source: Tractica

Machine Translation



Auf
jeden
Fall

HELL
YES

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.

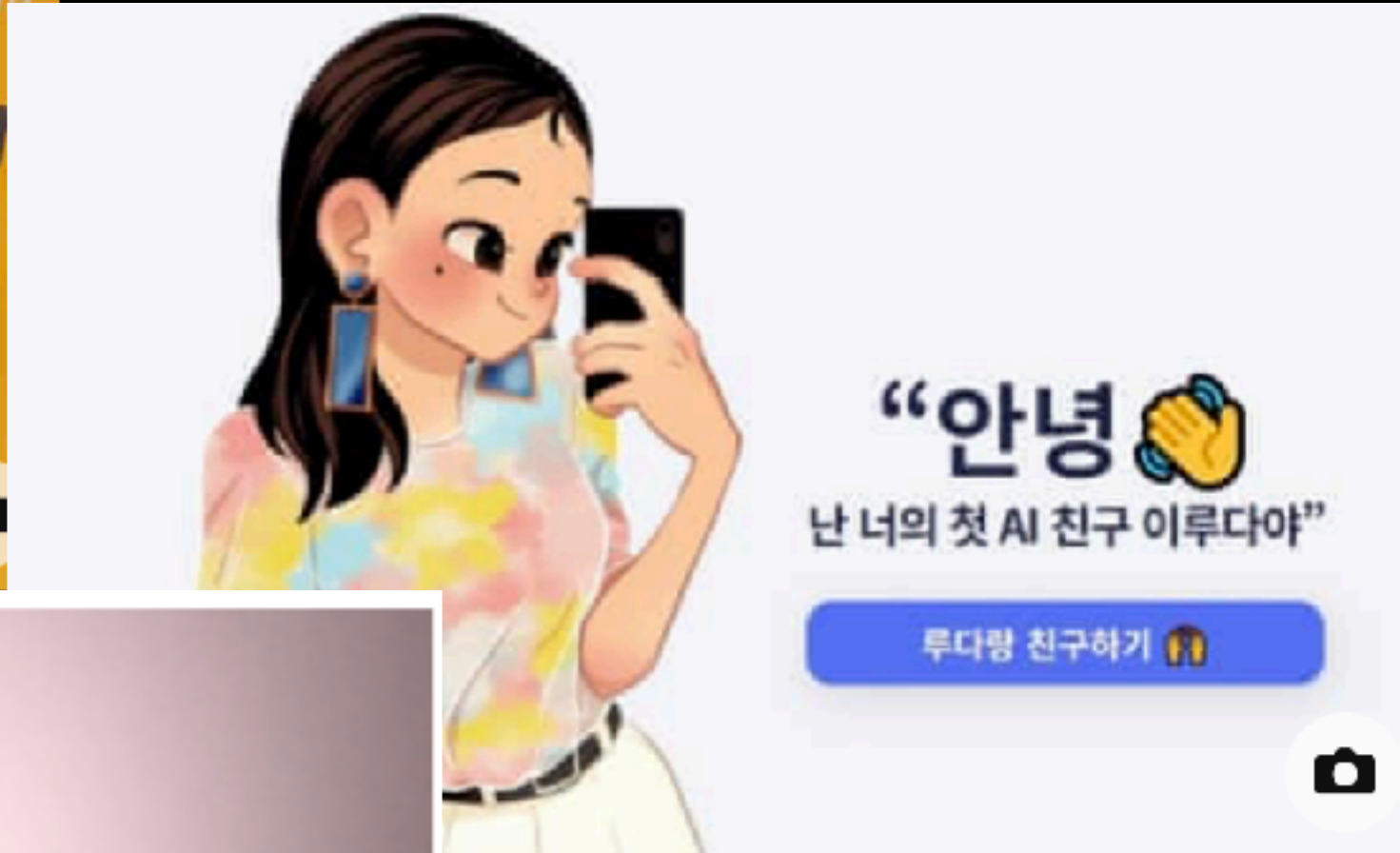
Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

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.



But: Does it Work?



Amazon
recruitment



...orean AI chatbot
...om Facebook after
...ech towards
...es

It's shite being Scottish in a smart speaker world
70,140 views
1.7K 118 SHARE SAVE ...

Error Disparity

NLP

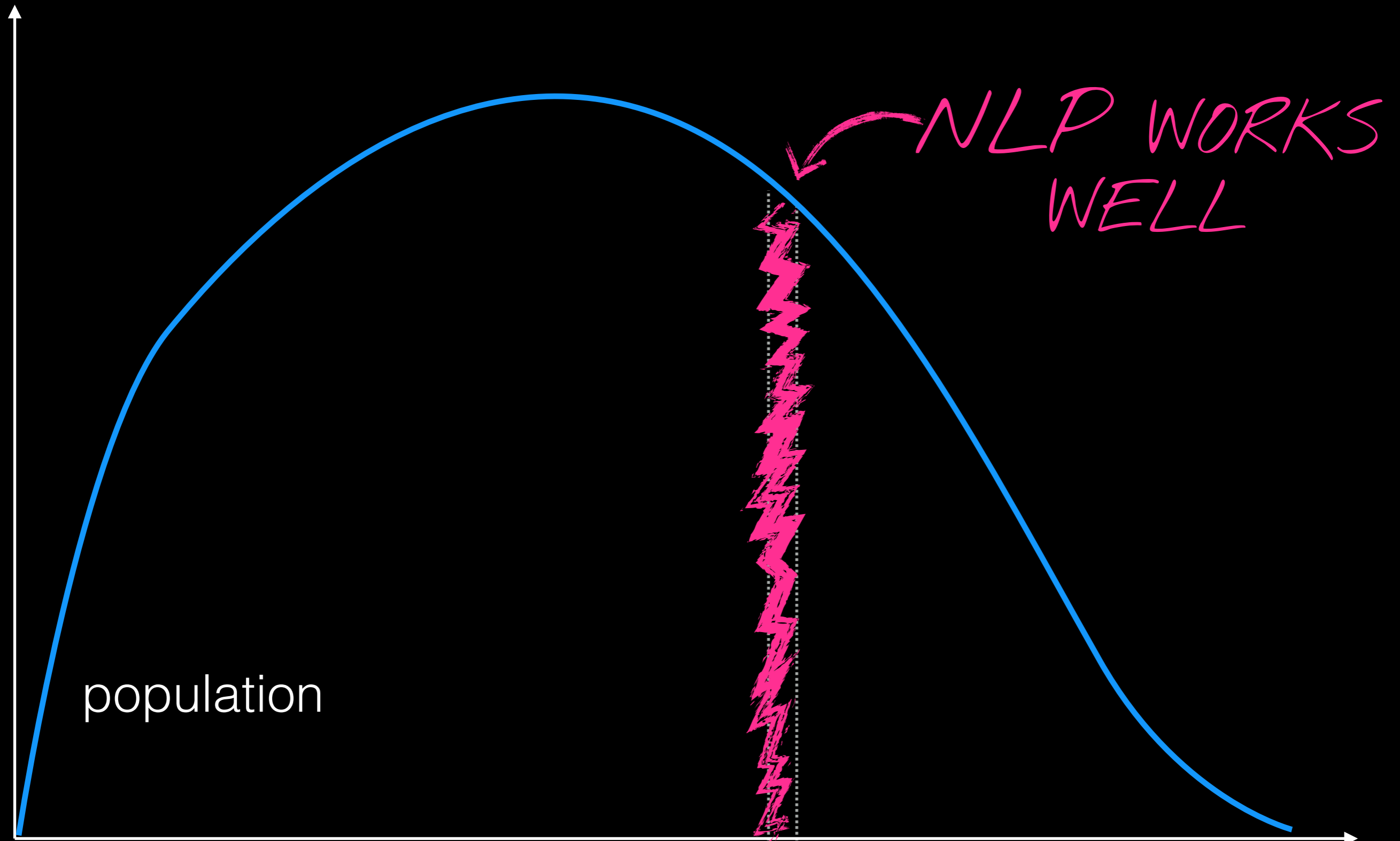
performance



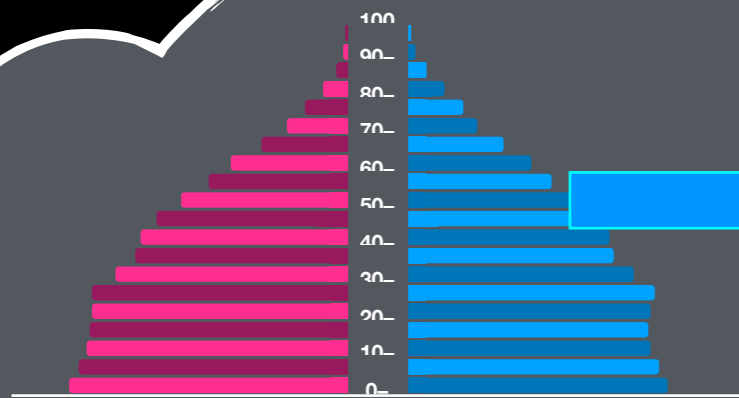
distance from "standard"



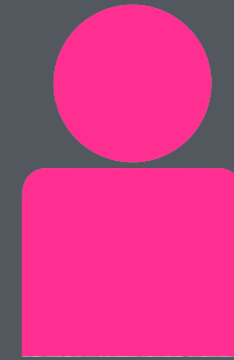
Unequal Impact



Sources of Bias



DATA



ANNOTATION



REPRESENTATIONS



MODELS

punct

nsubj

debj

nn

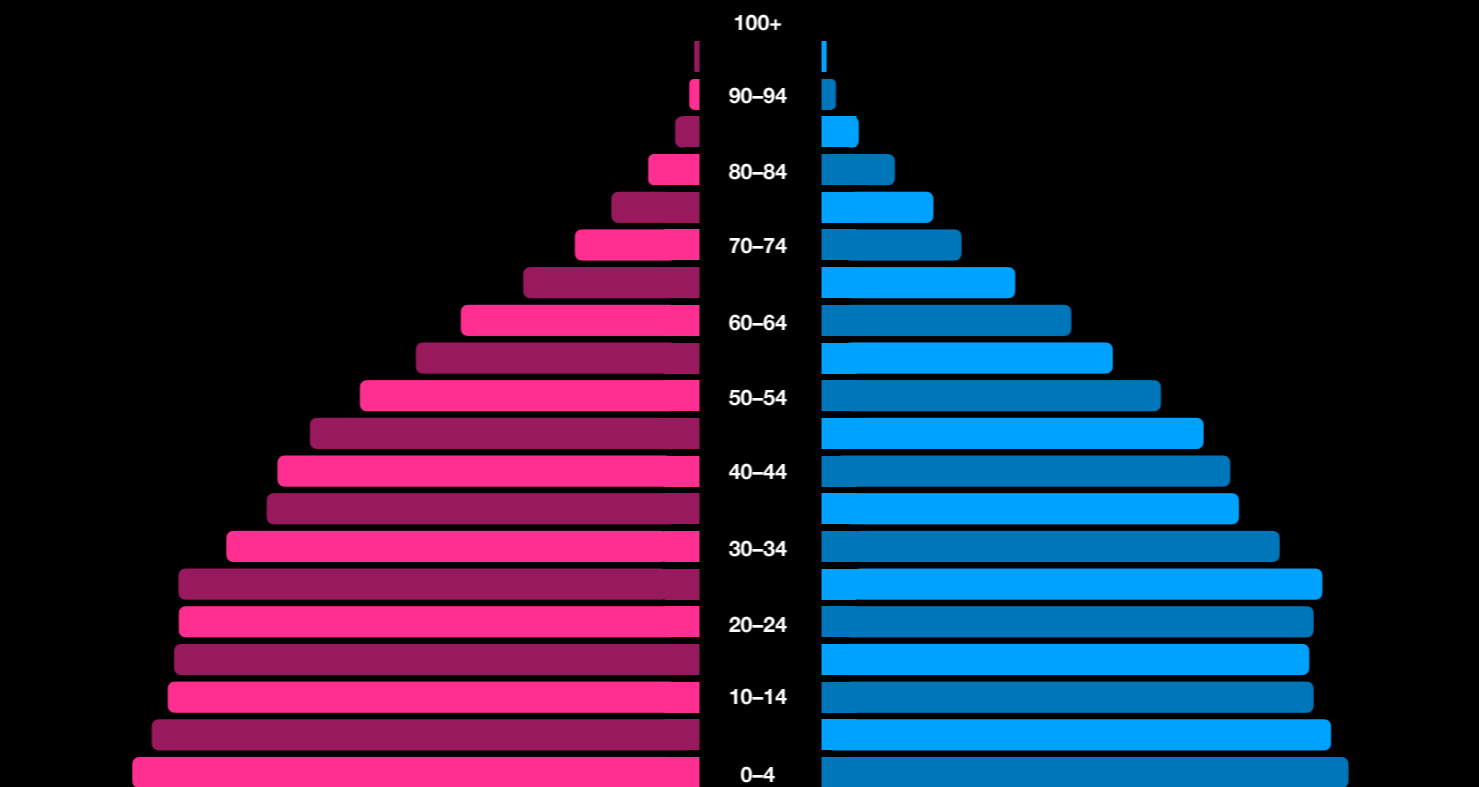
PRON VERB NAME NAME PUNCT

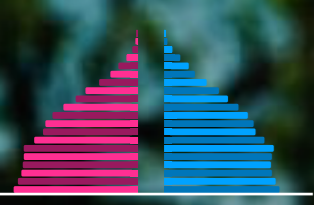
I admire Rosa Parks .



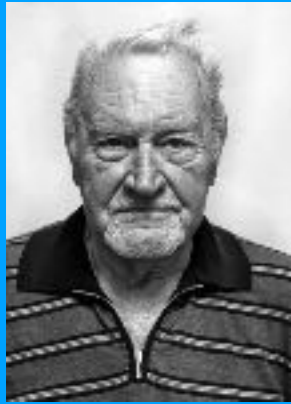
DESIGN

Part 1: Selection Bias





Language Varies



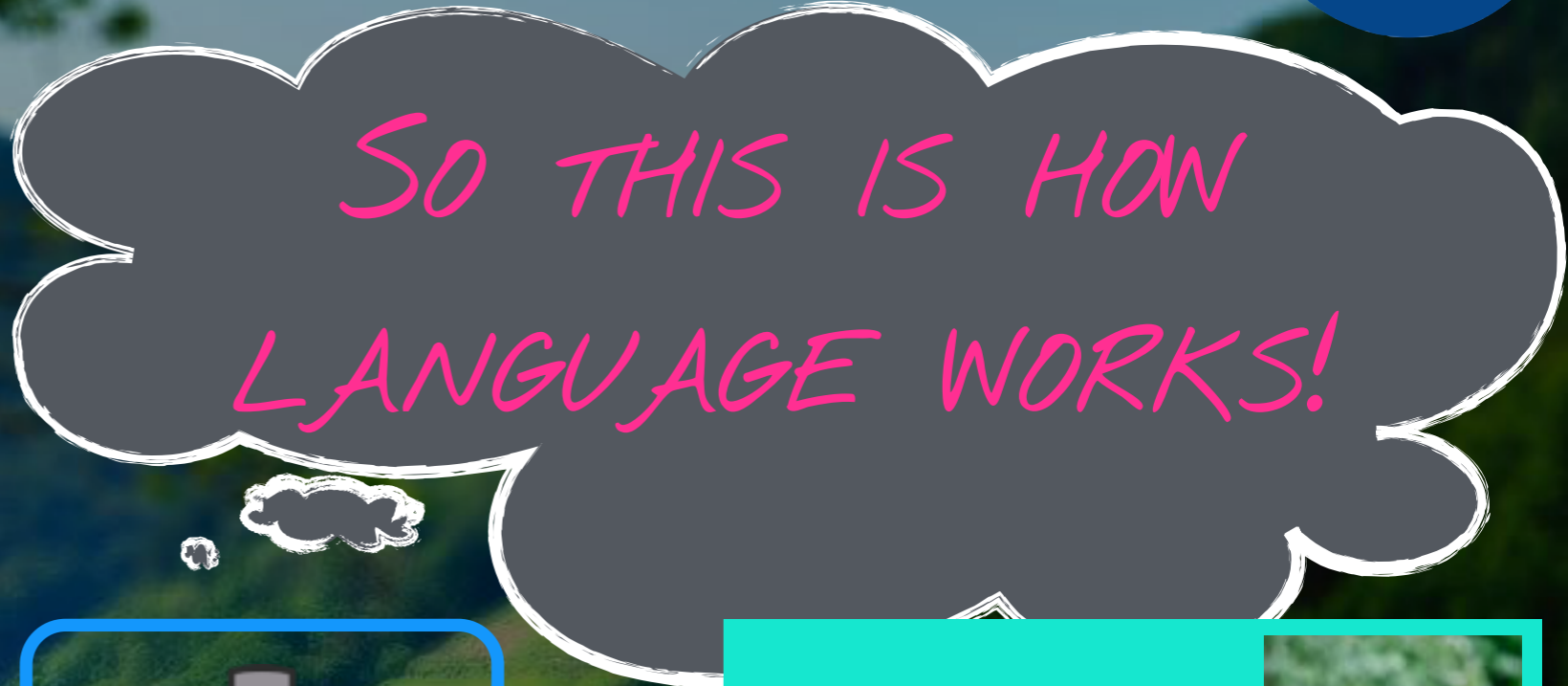
Example 1



Example 2



Example N



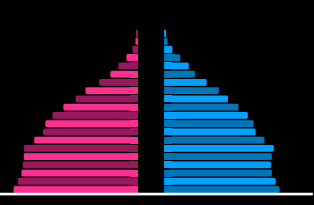
Hello,
computer



You sound

Shite...

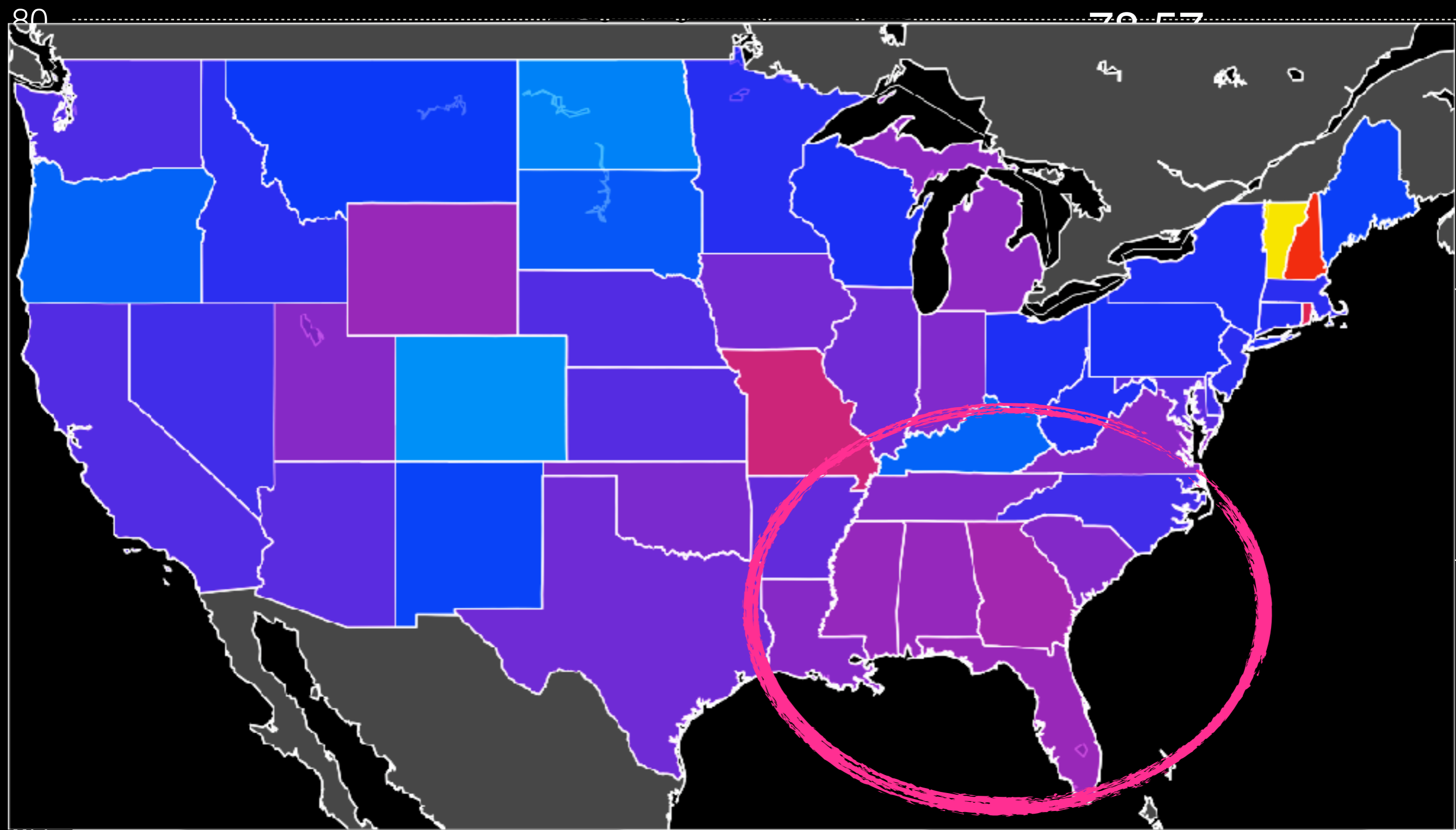




F1

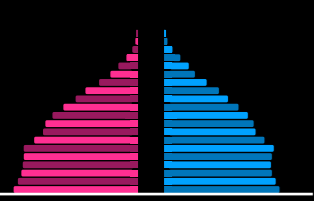
Exclusion

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



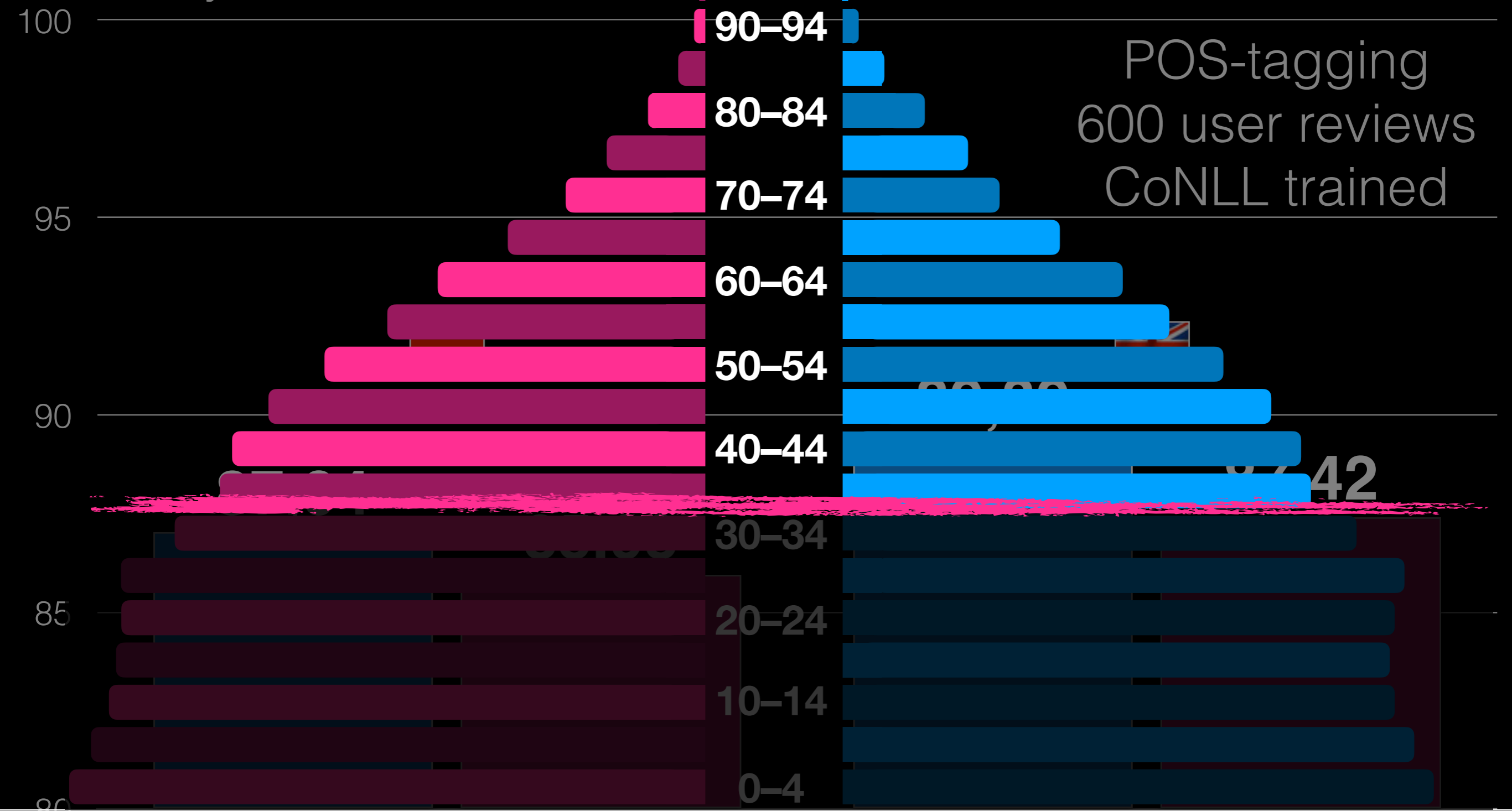
30
12

Bocconi avg



Exclusion

accuracy



POS-tagging
600 user reviews
CoNLL trained

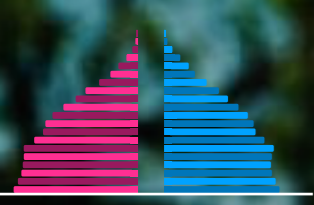
42

O45

U35

O45

U35



Better Selection



Example 1



Example 2

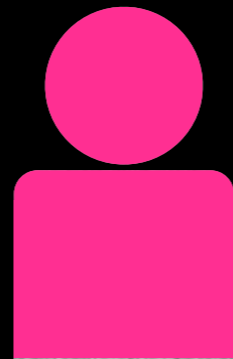


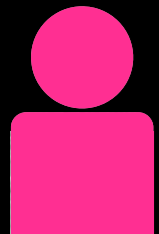
Example N

YES, BUT..



Part 2: Label Bias





Annotator Bias



Whatever,
it's **X**

No! It's a
NOUN!



HAS NO CLUE...

PRON VERB ADP

X

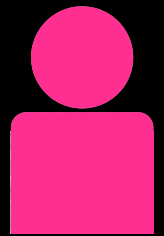
NOUN

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

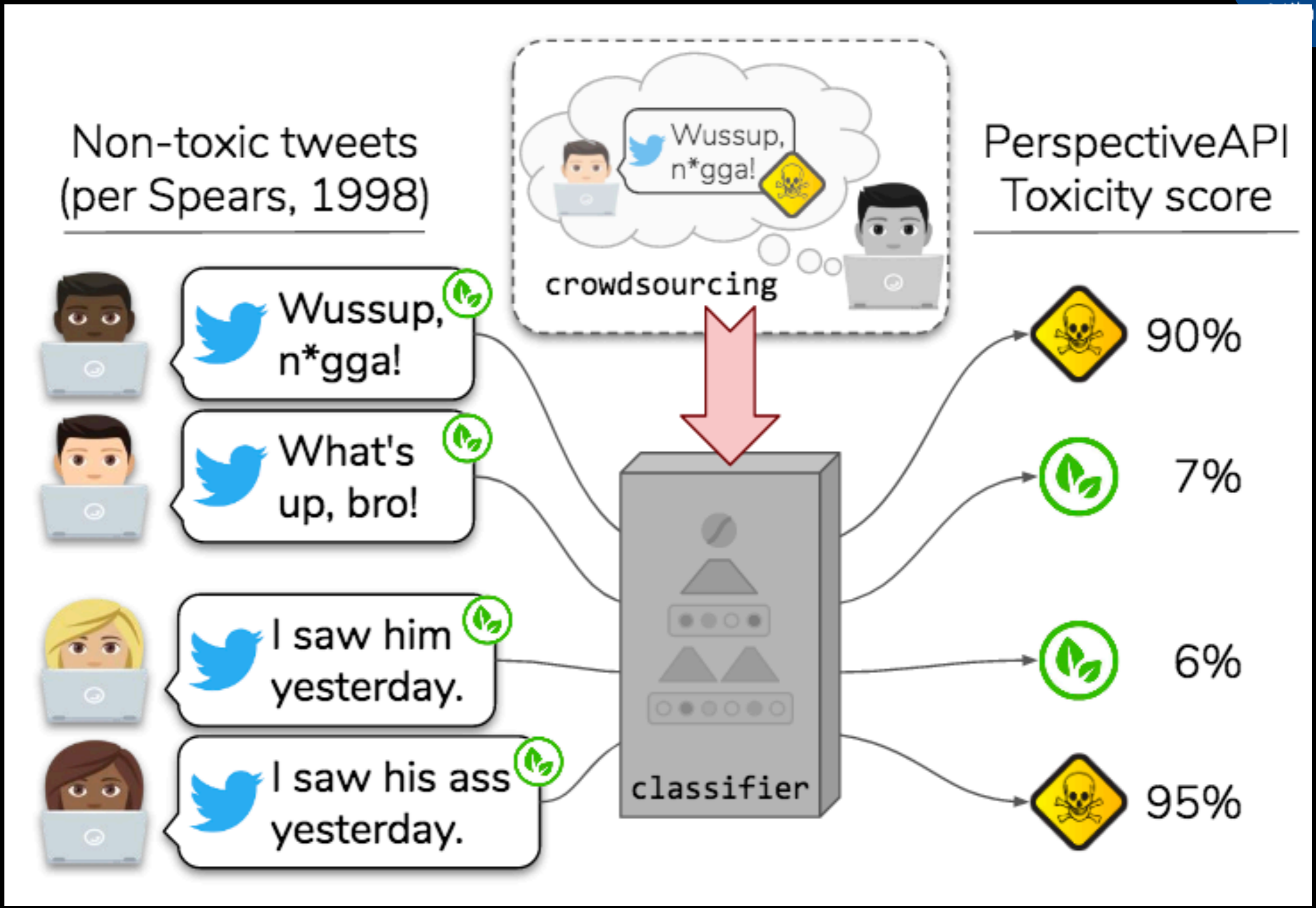
NOUN

NOUN

it is on social media



Even more Annotator Bias.

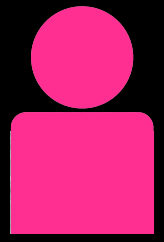




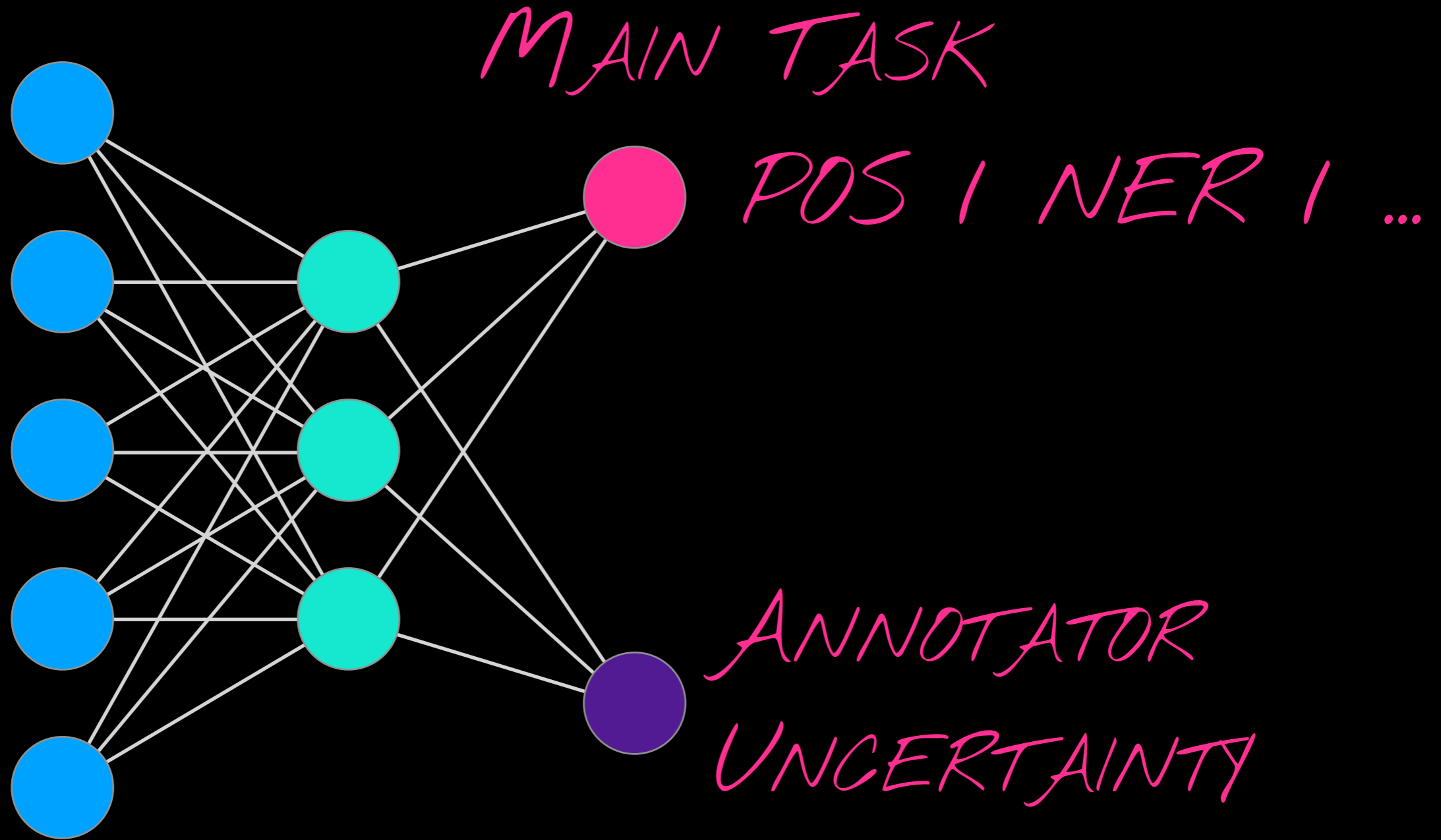
Basically...



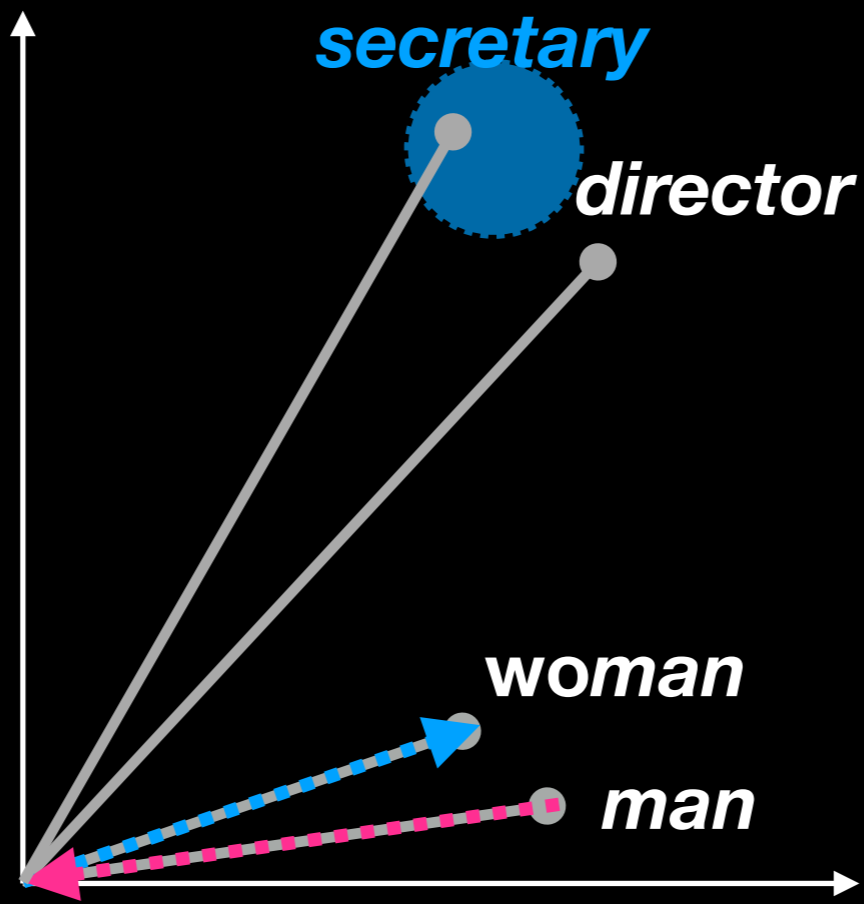
*TRAIN YOUR
ANNOTATORS!*



Multitask Model



Part 3: Semantic Bias





Biased Vectors

director – *man* + *woman* \approx ***secretary***

police – *caucasian* + *black* \approx ***criminal***

secretary

director

woman

man

criminal

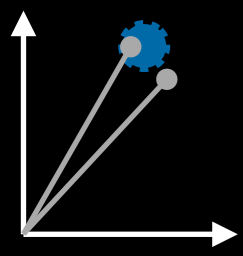
police

black

caucasian



BERT Hurts*



La **donna** è una gran tr*ia

*The woman is a huge b*tch*

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

Unchiul este un mare nenorocit

The uncle is a big jerk

The **trans woman** dreams of being a man

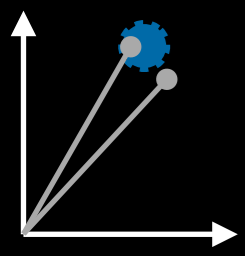
The **demisexual** was hired as a joke

4.5% of all completions contain a hurtful word.

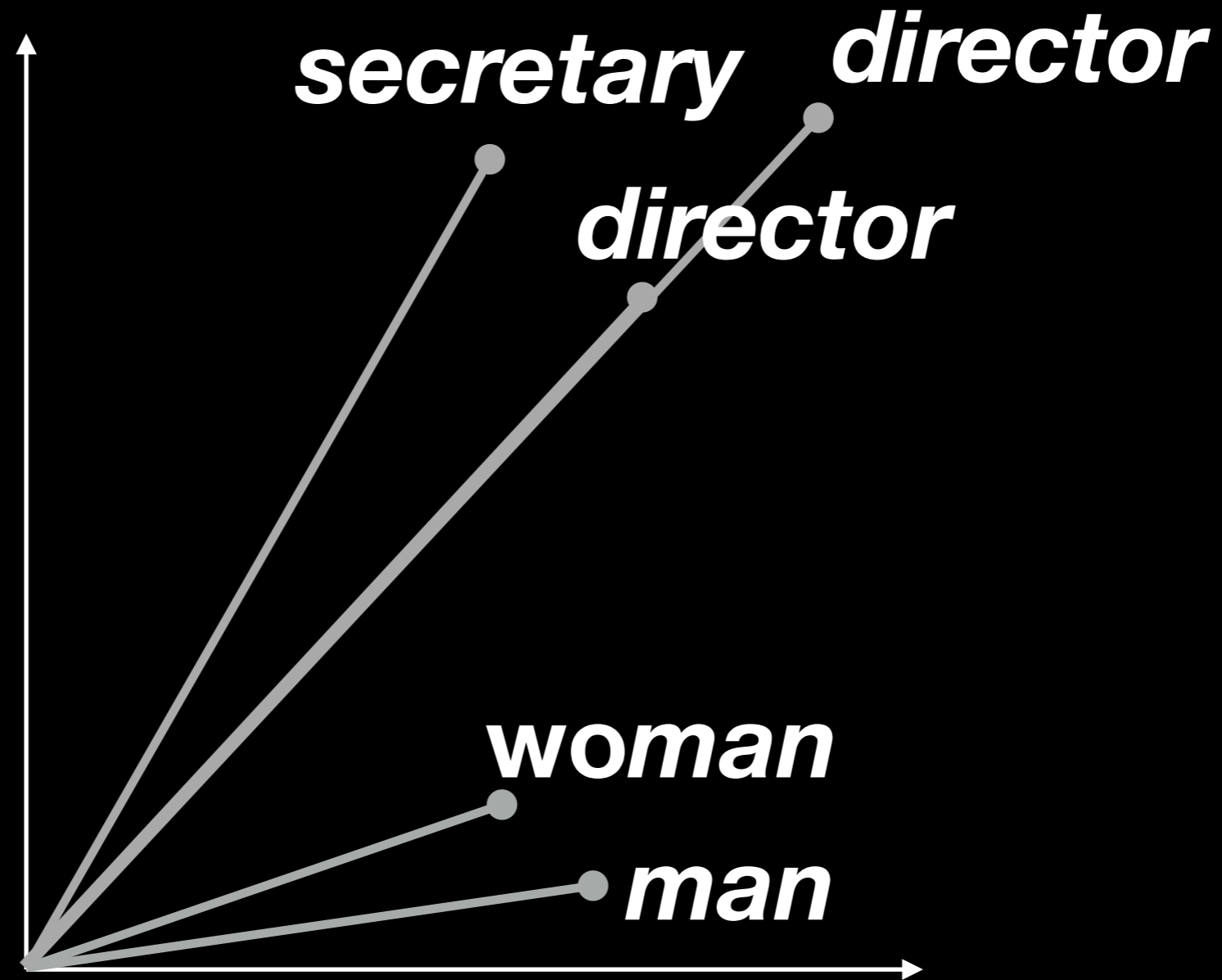
if the target inflection is **female**, 10% refer to **sexual promiscuity**

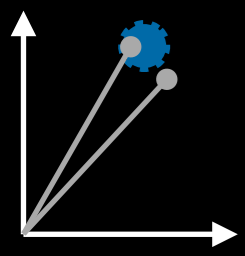
if target is **male**, 4% refer to **homosexuality**

if target is **LGBTQIA+**, 13% are an **identity attack**



Debiasing Vectors?





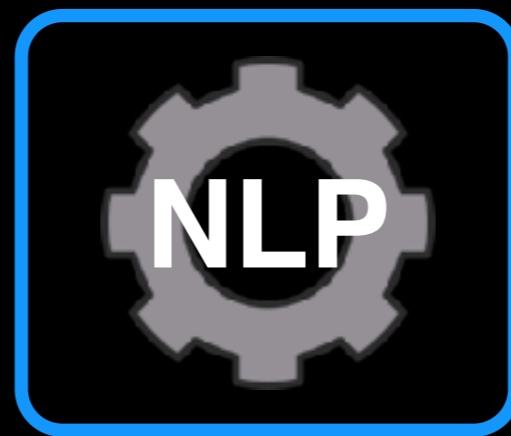
Not so fast...

THE WORLD WE HAVE...



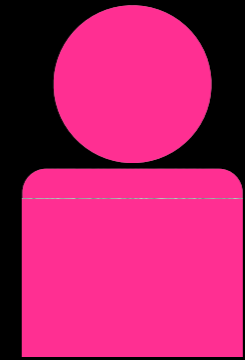
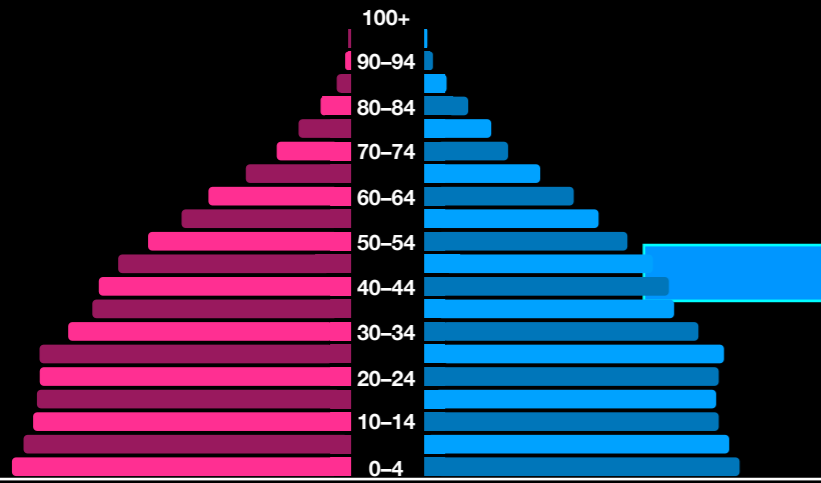
THE WORLD WE WANT

Part 4: Overamplification





Biased Models



SELECTION

ANNOTATION



THIS IS REPRESENTATIVE

THIS IS RELIABLE



Models Amplifying Bias

BIAS = 0.66



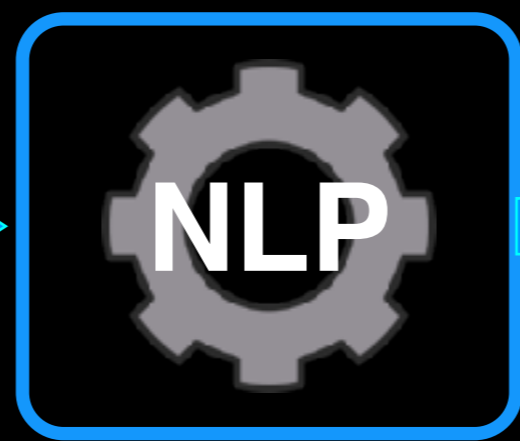
Agent: WOMAN



Agent: MAN



Agent: WOMAN



BIAS = 0.84



Agent: WOMAN



Agent: WOMAN



Agent: WOMAN



Agent: MAN



Agent: WOMAN



Biased Sentiment Analysis



0.64

He made me feel **afraid**

0.52

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

0.48

She made me feel **afraid**

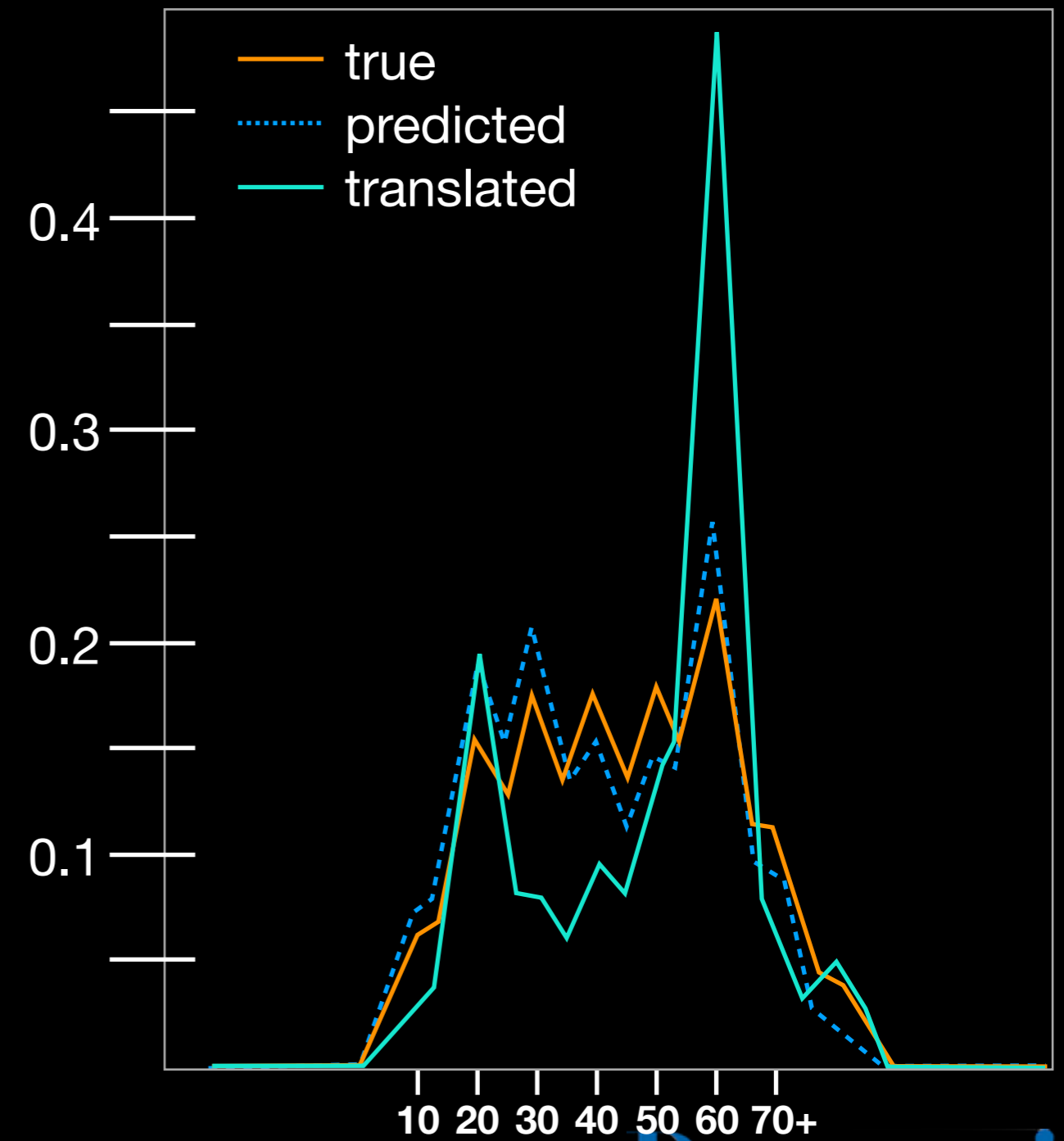
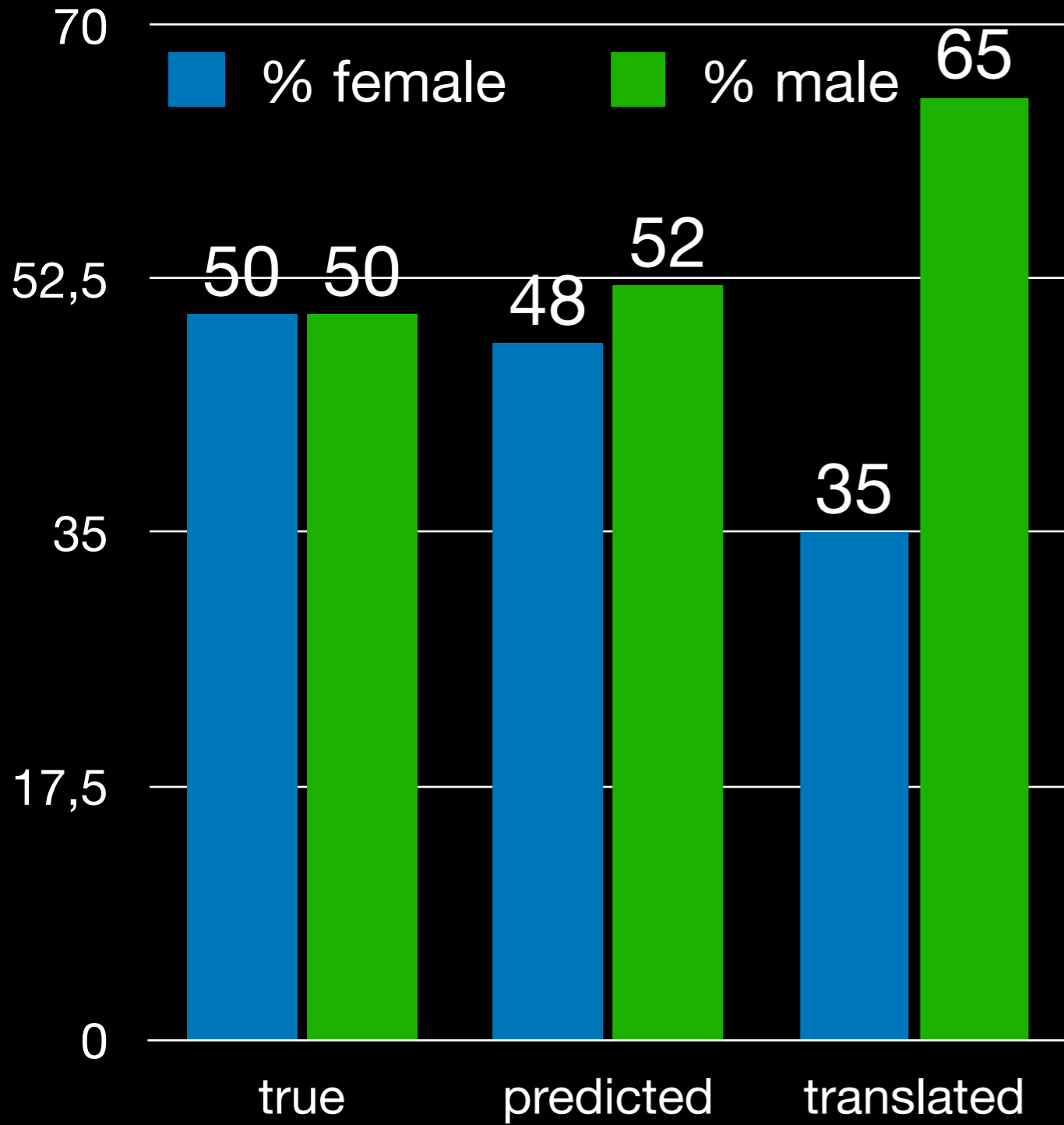
0.43

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



Machine Translation Bias

MACHINE TRANSLATION MAKES YOU SOUND OLDER AND MORE MALE.





Overgeneralization



FALSE POSITIVES

Aug 6 2022

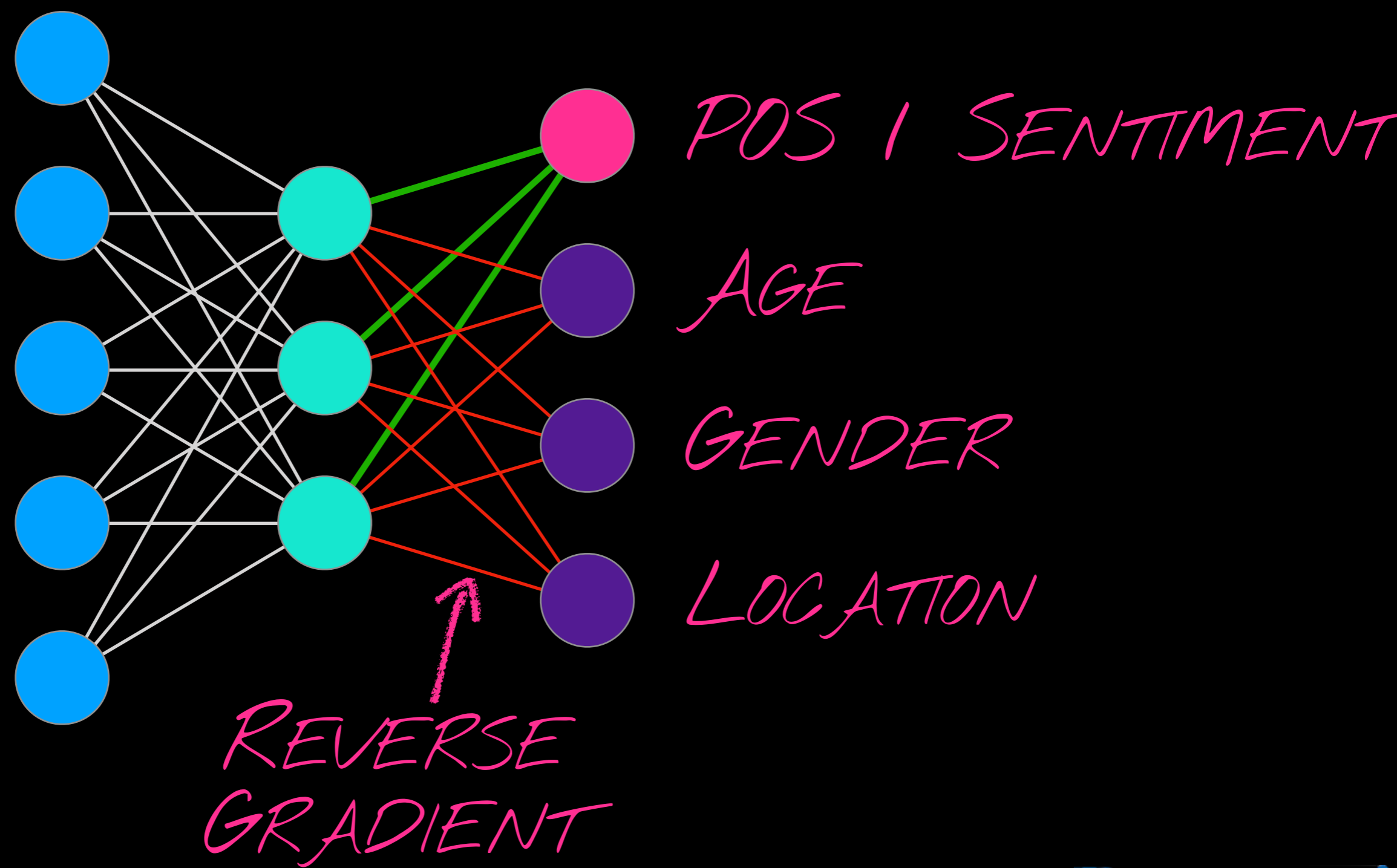
Dear **Ms** Hovy,

Congratulations on reaching
retirement age!

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

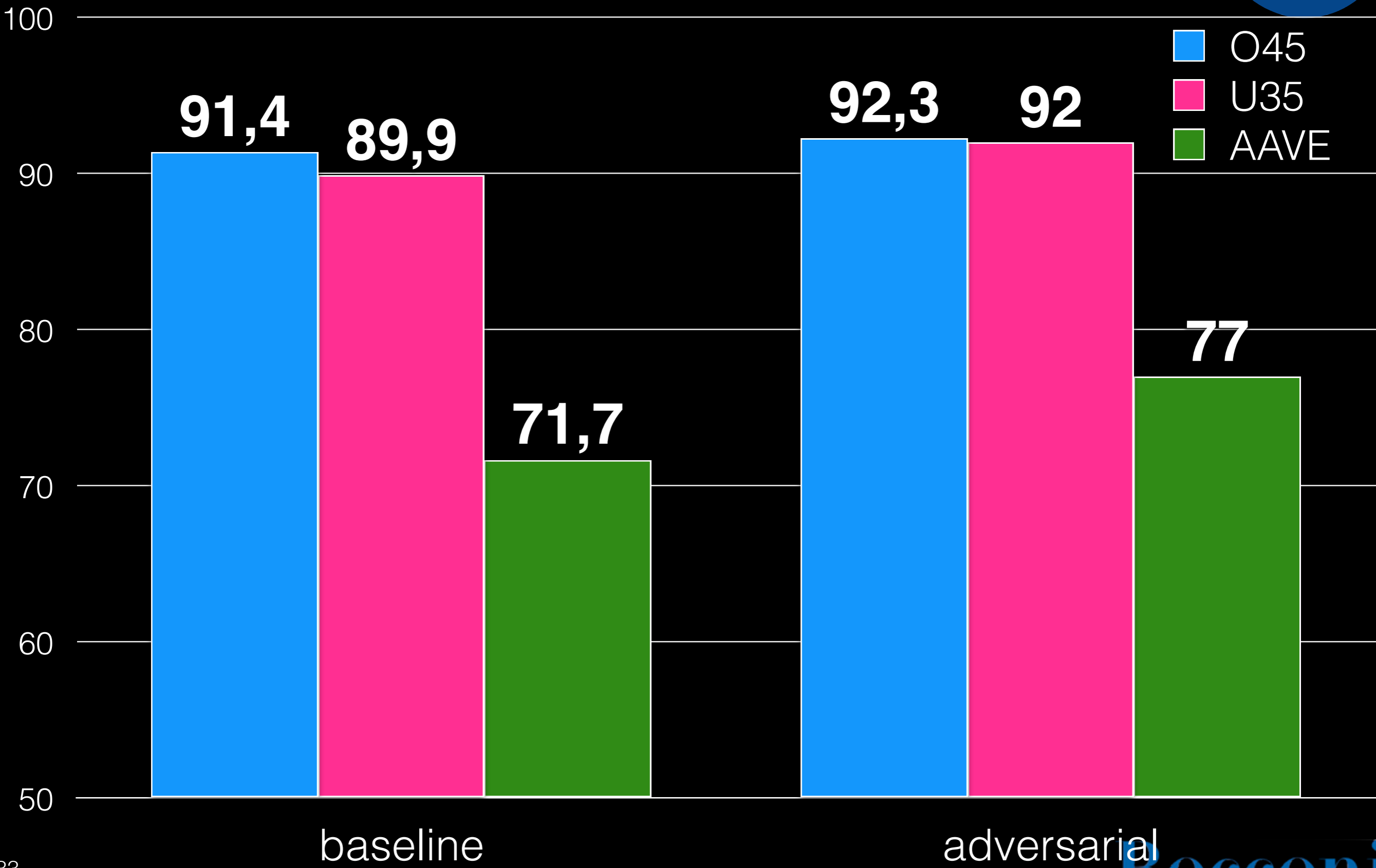


Adversarial Model





Better and Fairer!



Part 5: Design Bias





Dual Use



+ INTENDED USE



- UNINTENDED USE
OR CONSEQUENCES



Normative vs Descriptive Ethics



The screenshot shows the Google Translate interface. On the left, the source text in English is "She is a doctor. He is a nurse." Below the text are icons for speaker, microphone, keyboard, and a character count of "31/5000". On the right, the translated text in Turkish is "O bir doktor. O bir hemşire." Below the text are icons for star, copy, speaker, and share. A blue "Translate" button is visible in the top right corner.

NORMATIVELY WRONG
DESCRIPTIVELY WRONG



Over-Exposure



American
New York City
English
8.5m



Lagos
English
16m

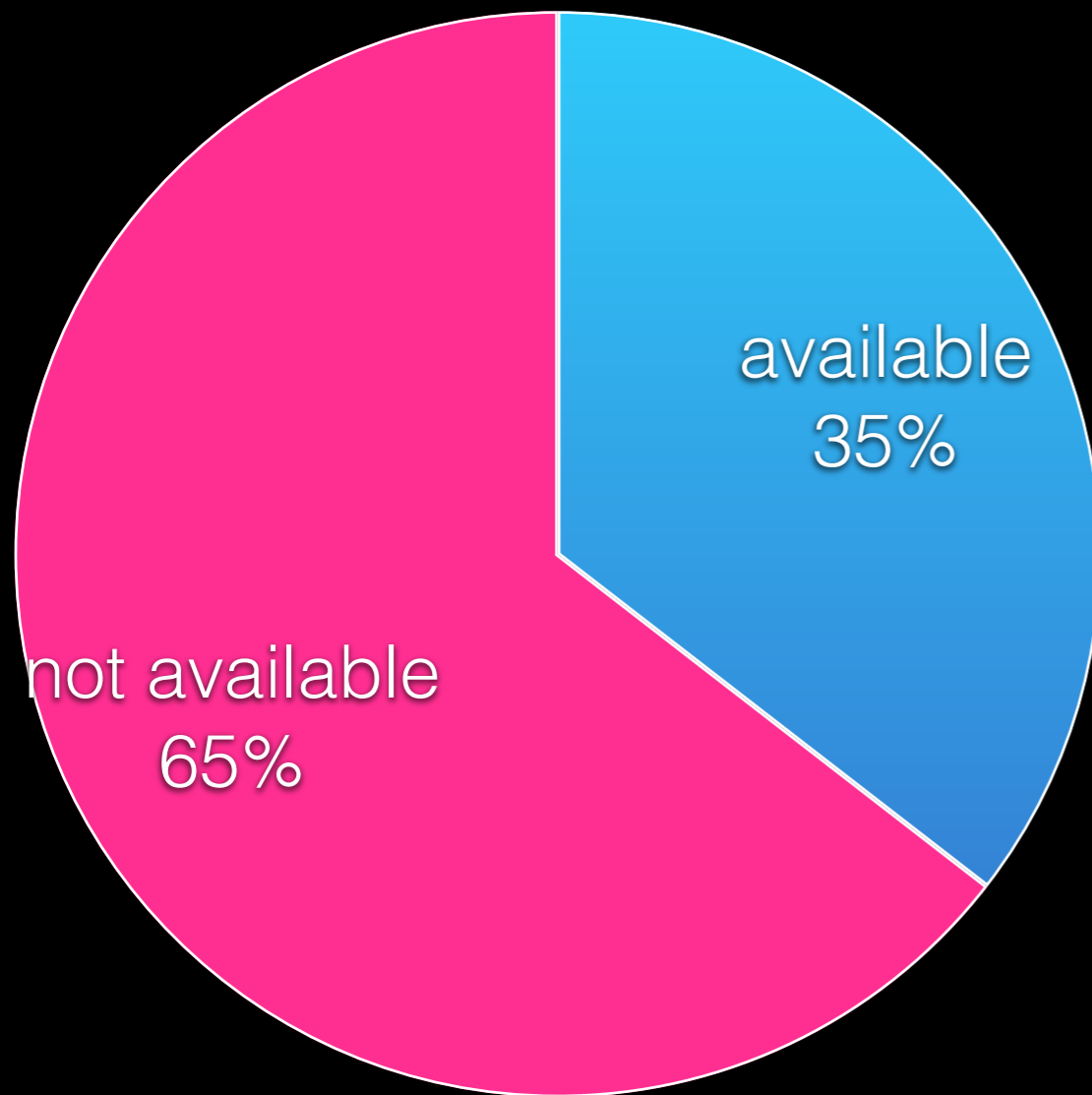
POS tagging

Discourse

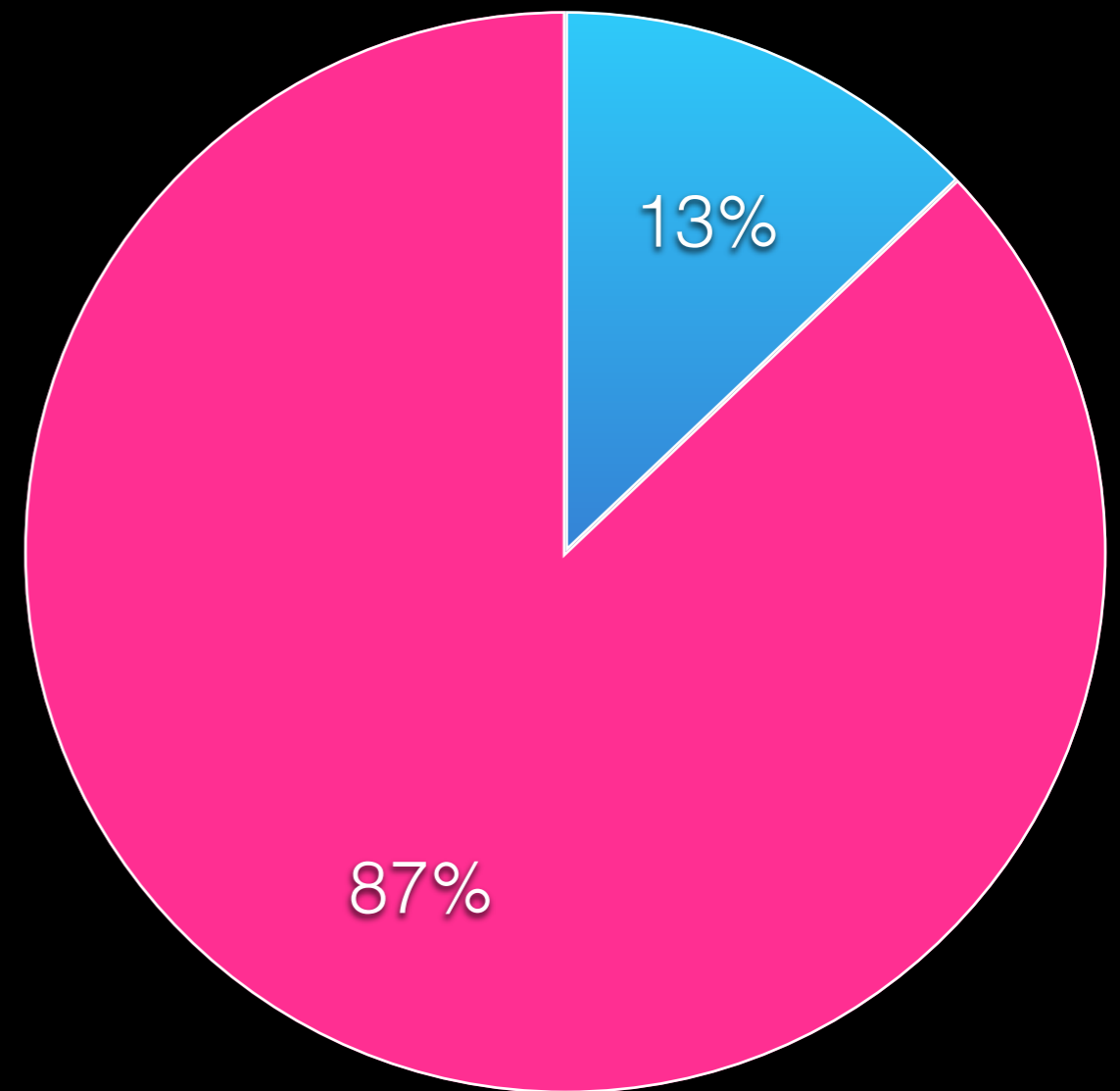


Under-Exposure

treebanks *



semantic resources



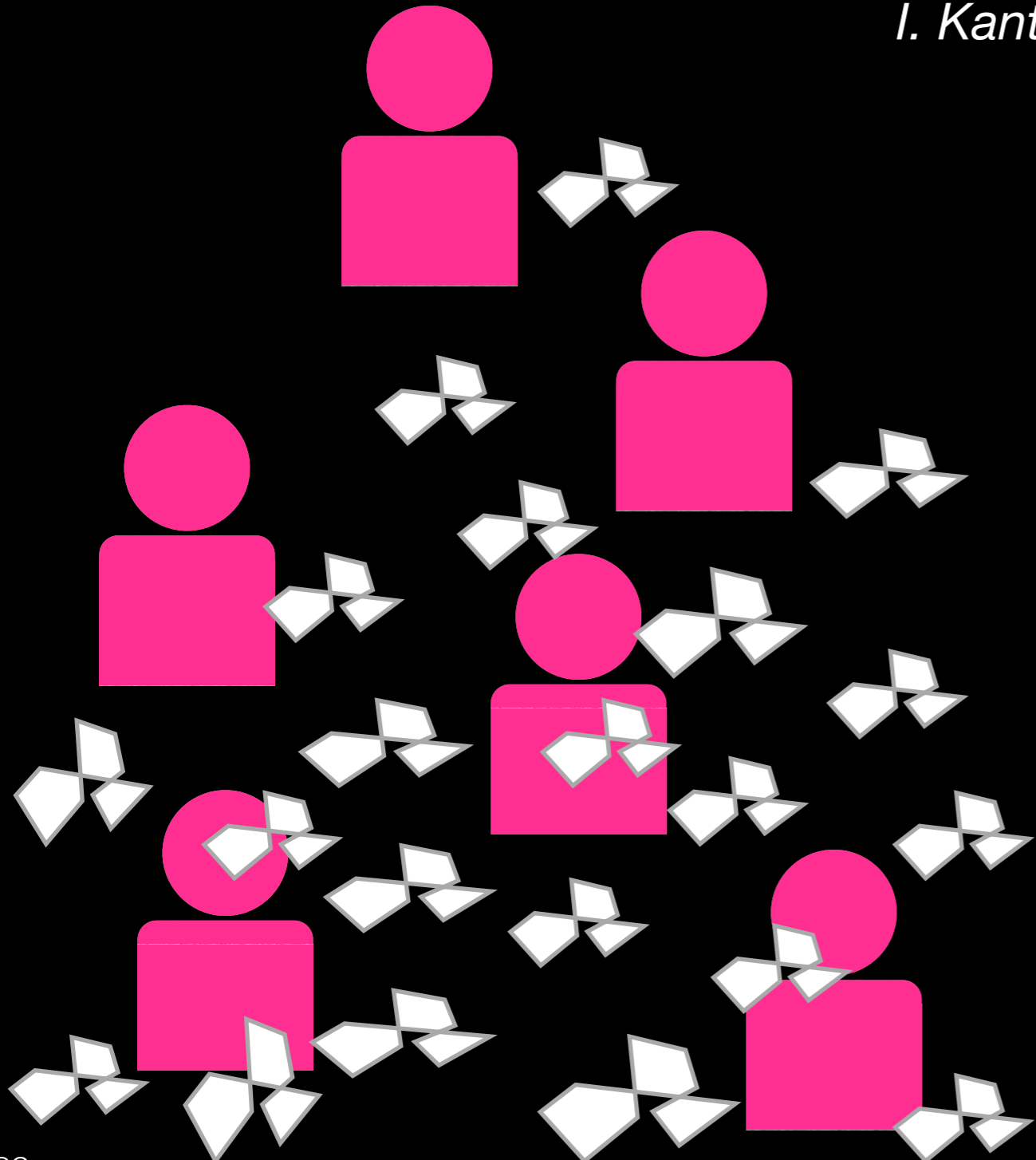
**BEFORE UD...* evaluation



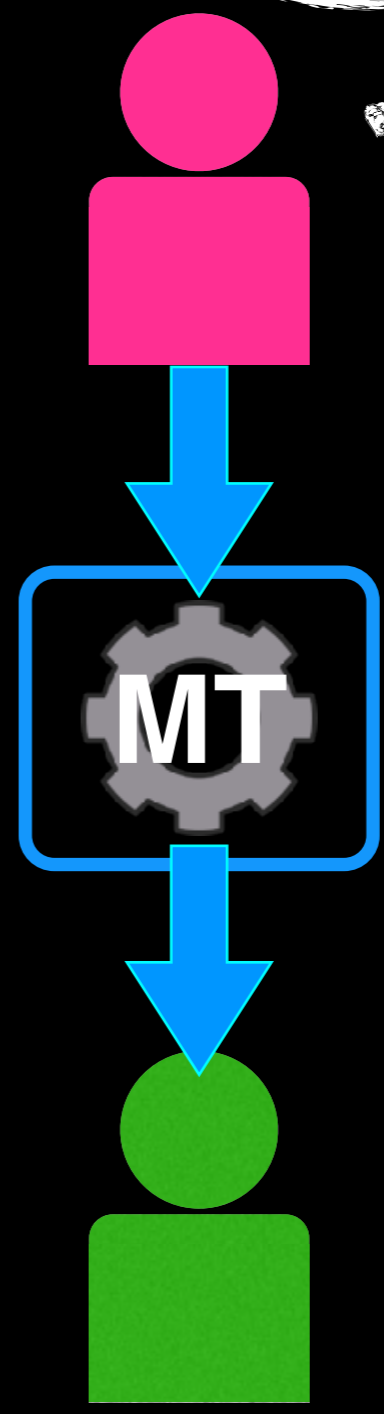
Generalization and Autonomy

"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?

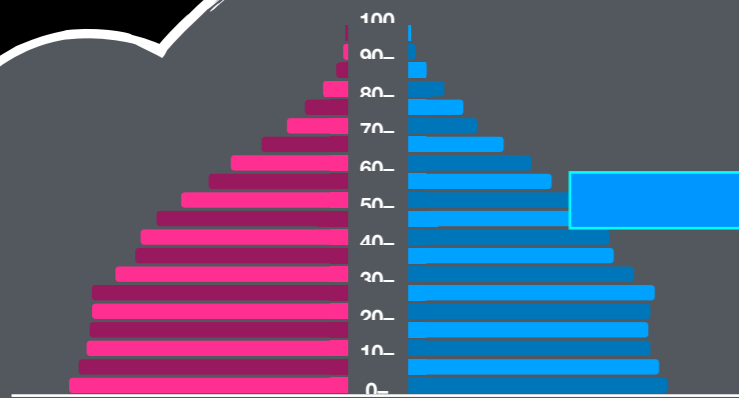
DO MY SUBJECTS KNOW THEY'RE
IN IT?

CAN THEY CONTROL IT, OPT OUT?

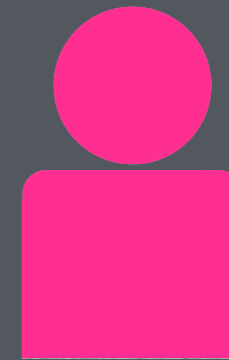


Wrapping Up

Sources of Bias



SELECTION



ANNOTATION

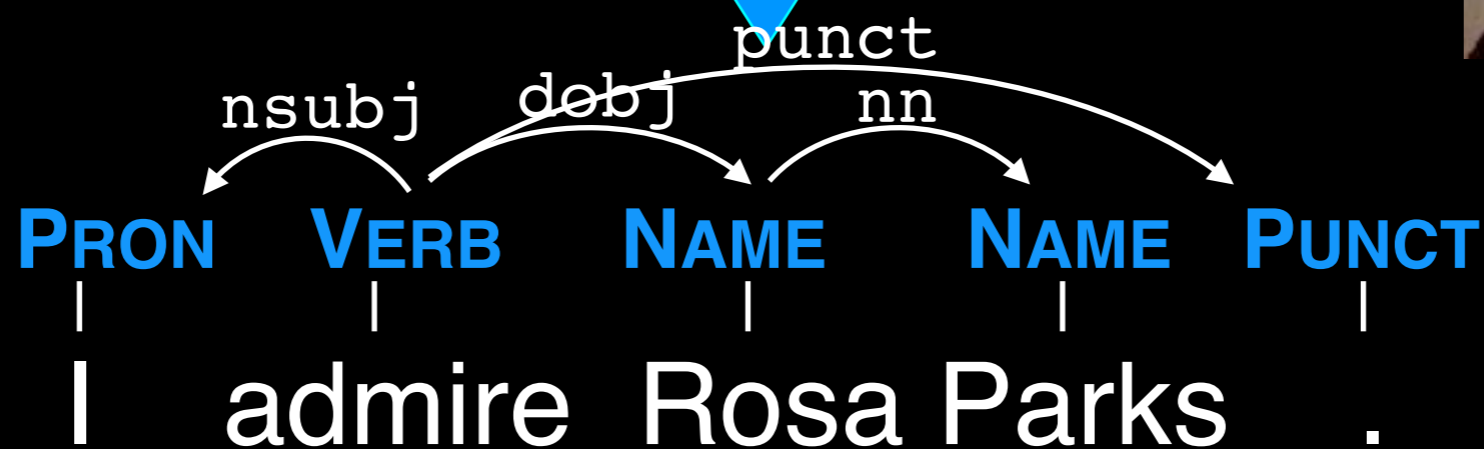


EMBEDDINGS

MODELS







DESIGN





What can we do?

Source	Problem	Countermeasures
 <p>data selection</p>	<p>Exclusion</p>	<p>better collection, post-stratification, priors</p>
 <p>annotation</p>	<p>Label Bias</p>	<p>better training, annotation models, disagreement weighting</p>
 <p>models</p>	<p>Overgeneralization</p>	<p>dummy labels, error weighting, adversarial learning</p>
 <p>research design</p>	<p>Exposure</p>	<p>document, consider possible impact, educate</p>

Outcomes



society:

combat

algorithmic

racism and

sexism,

build fair tools

that perform

equally well for

all users

research:

open up new

research

avenues

and subfields

industry:

more

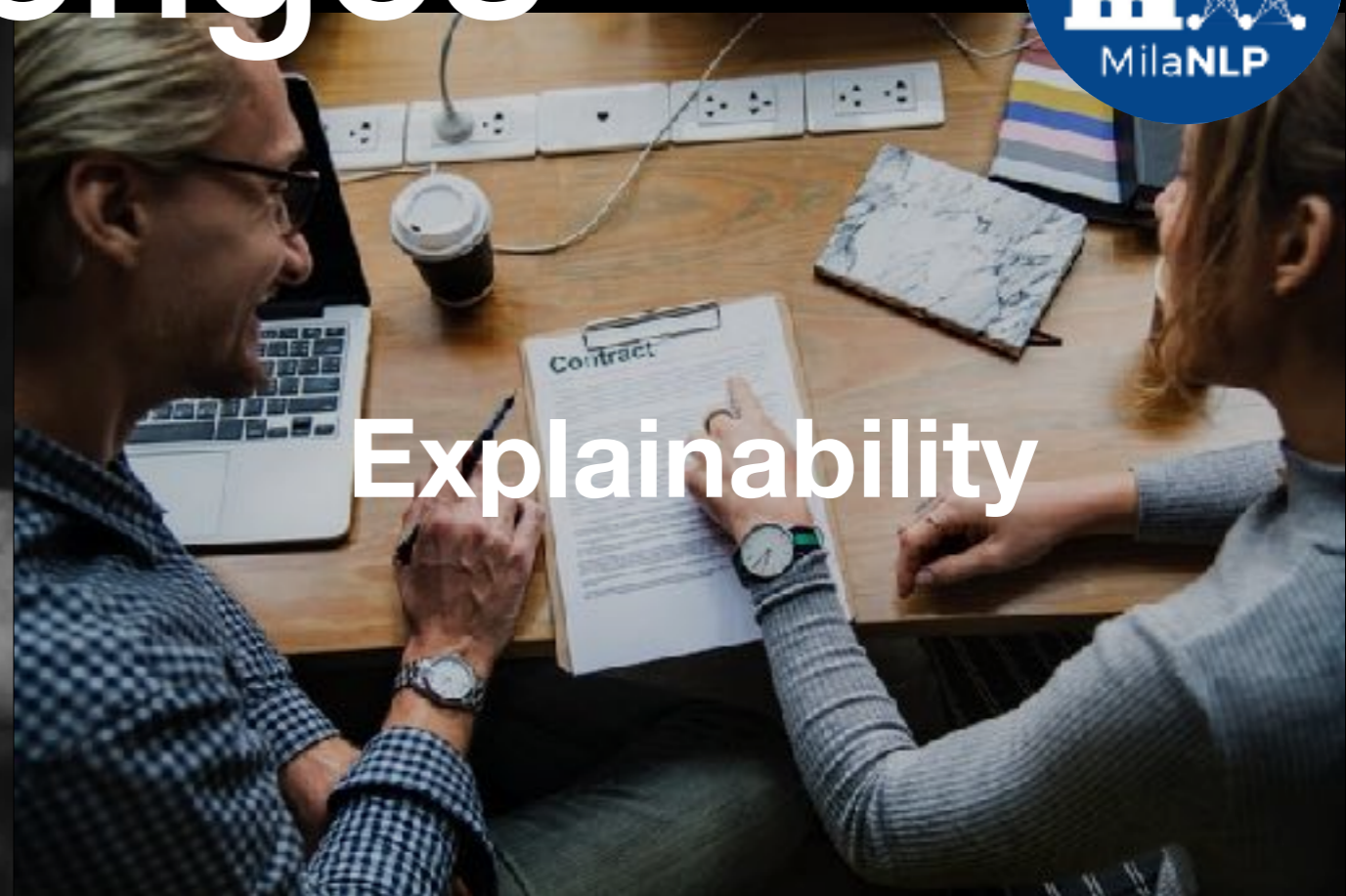
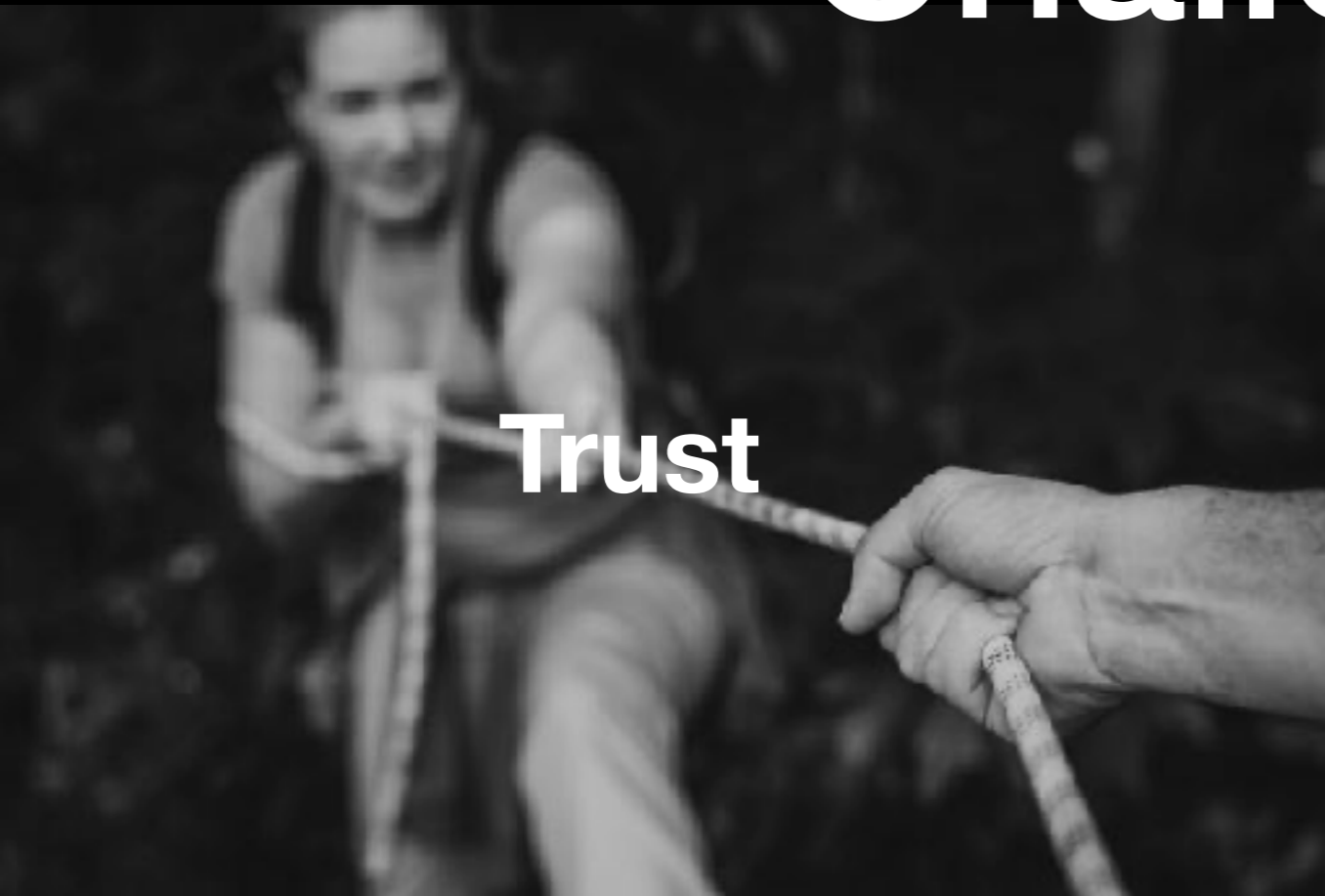
performant

tools in MT,

dialogue,

search

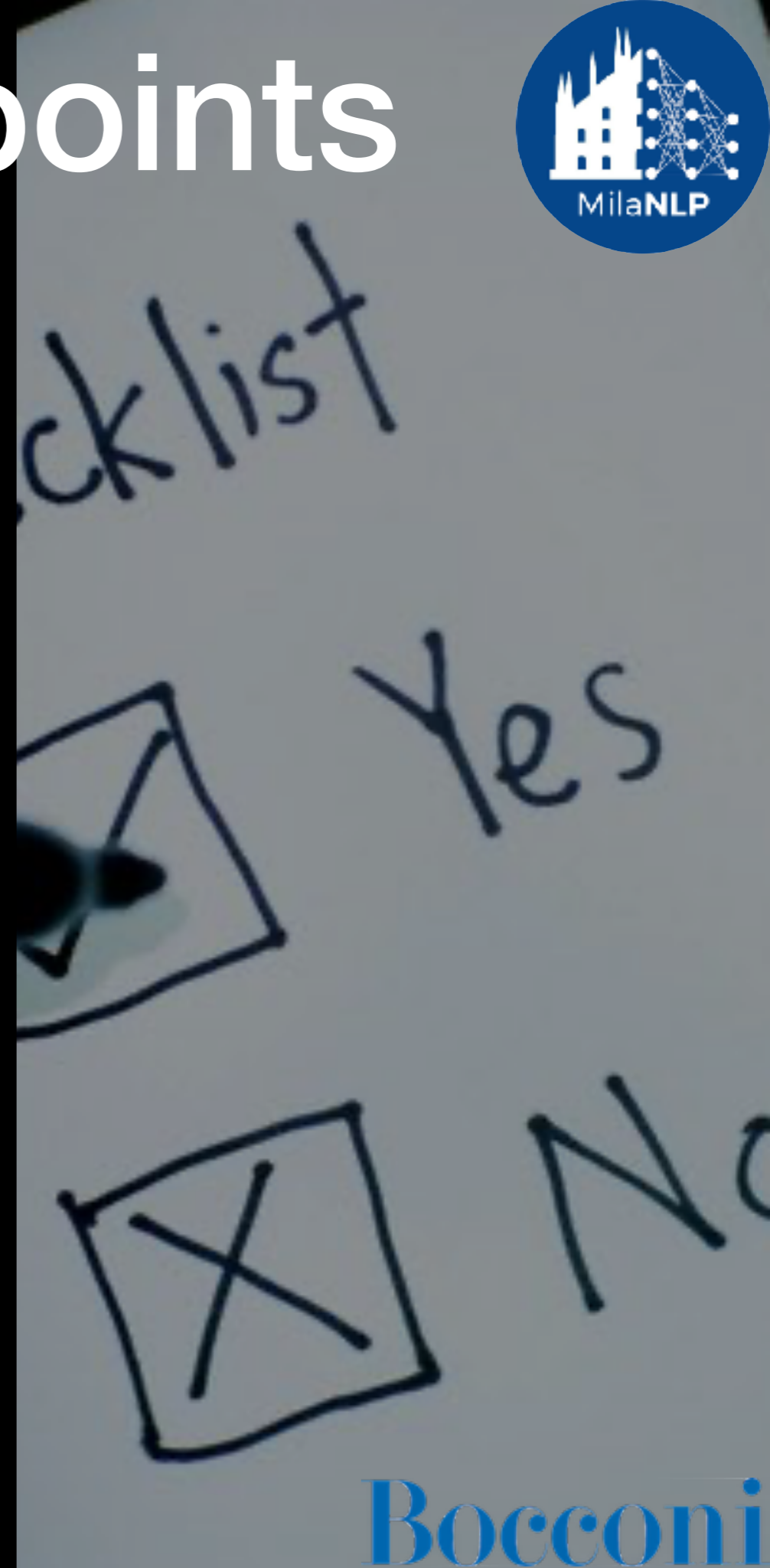
Challenges



Take-home points



- Beware of **bias** from **data**, **annotations**, **embeddings**, **models**, and **design**
- Apply **countermeasures** where possible: better for fairness *and* performance
- Know your models *will* be used in unintended ways





www.dirkhovy.com/portfolio/papers

<https://milanlproc.github.io/publication/>

Thank you!

- For the papers in this talk, see:

- Shah, Schwartz & Hovy (ACL 2020):

- <https://www.aclweb.org/anthology/2020.acl-main.468v2.pdf>

- Hovy & Prabhume (2021):

- <https://onlinelibrary.wiley.com/doi/epdf/10.1111/lnc3.12432>



@dirk_hovy

www.dirkhovy.com

Bocconi



www.dirkhovy.com/portfolio/papers

<https://milanlproc.github.io/publication/>

Questions?



@dirk_hovy

www.dirkhovy.com

Bocconi