



European Research Council  
Established by the European Commission



# More than Words: Integrating Social Factors into NLP

**Dirk Hovy**  
Bocconi University  
Milan, Italy

`dirk.hovy@unibocconi.it`  
`http://www.dirkhovy.com`

 `@dirk_hovy`

# Choices, Choices



I am **excited** to be here

So **pumped** to be here



# OK, Computer...



Goal: to make computers understand language like we do!

This talk: taking stock and looking ahead

# Since then...



Auf  
jeden  
Fall

HELL  
YES

# Since then...



*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.

# NLP is booming



6 Source: Tractica



# But: Does it Work?



*MEH...*

I'm trying to figure out if these devices will understand a single word I say.

It's shite being Scottish in a smart speaker world

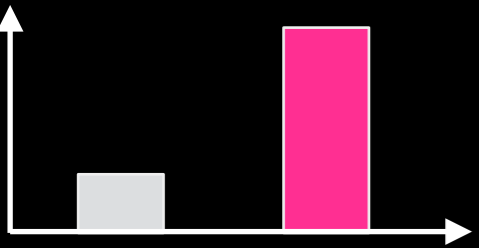
70,140 views

1.7K 118 SHARE SAVE ...

# Current Limitations



- creates demographic biases that disadvantage large segments of society



- negatively affects performance of existing NLP tools

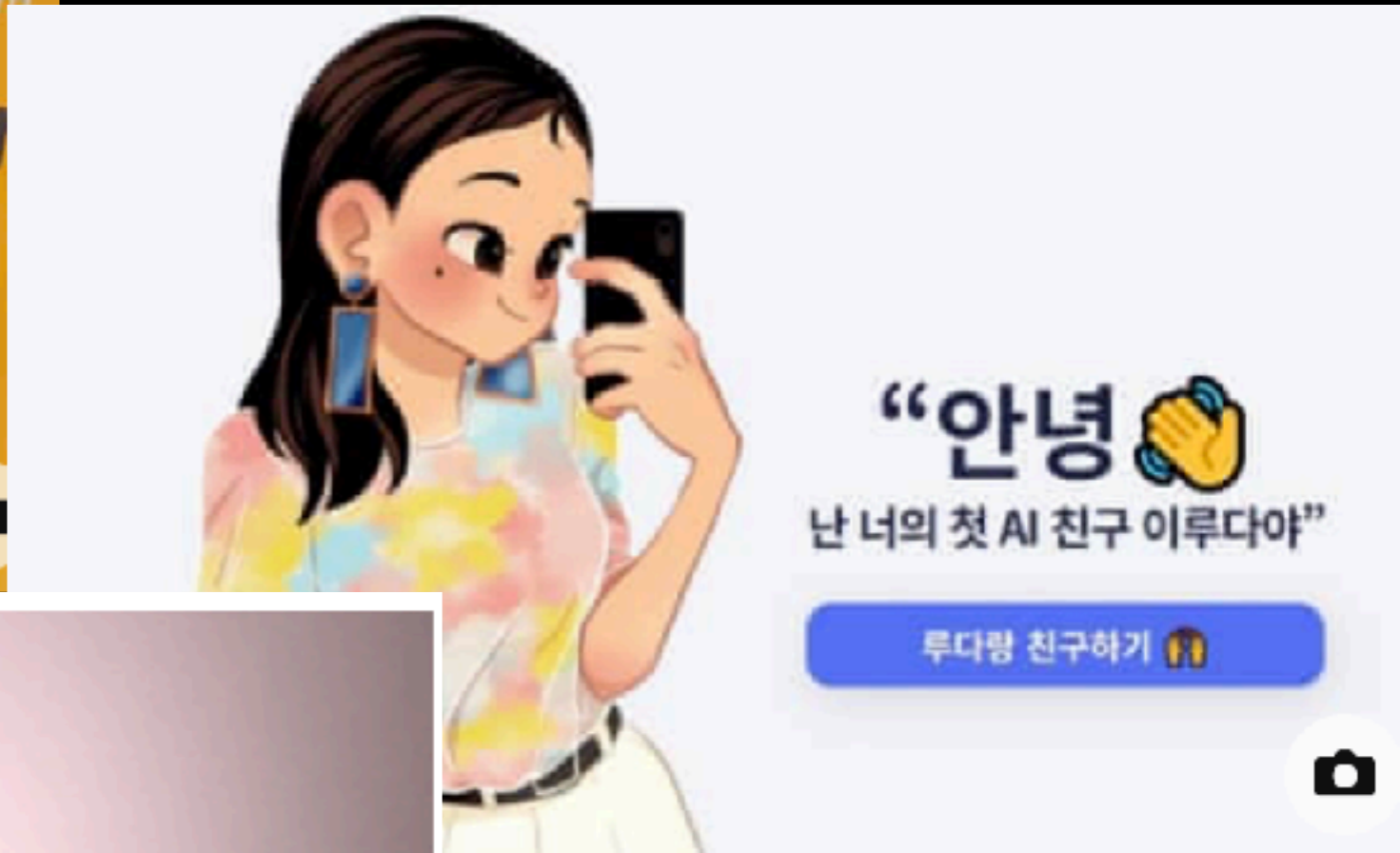


- severely limits novel applications to information-based tasks





# Biases



Amazon  
recruitment



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

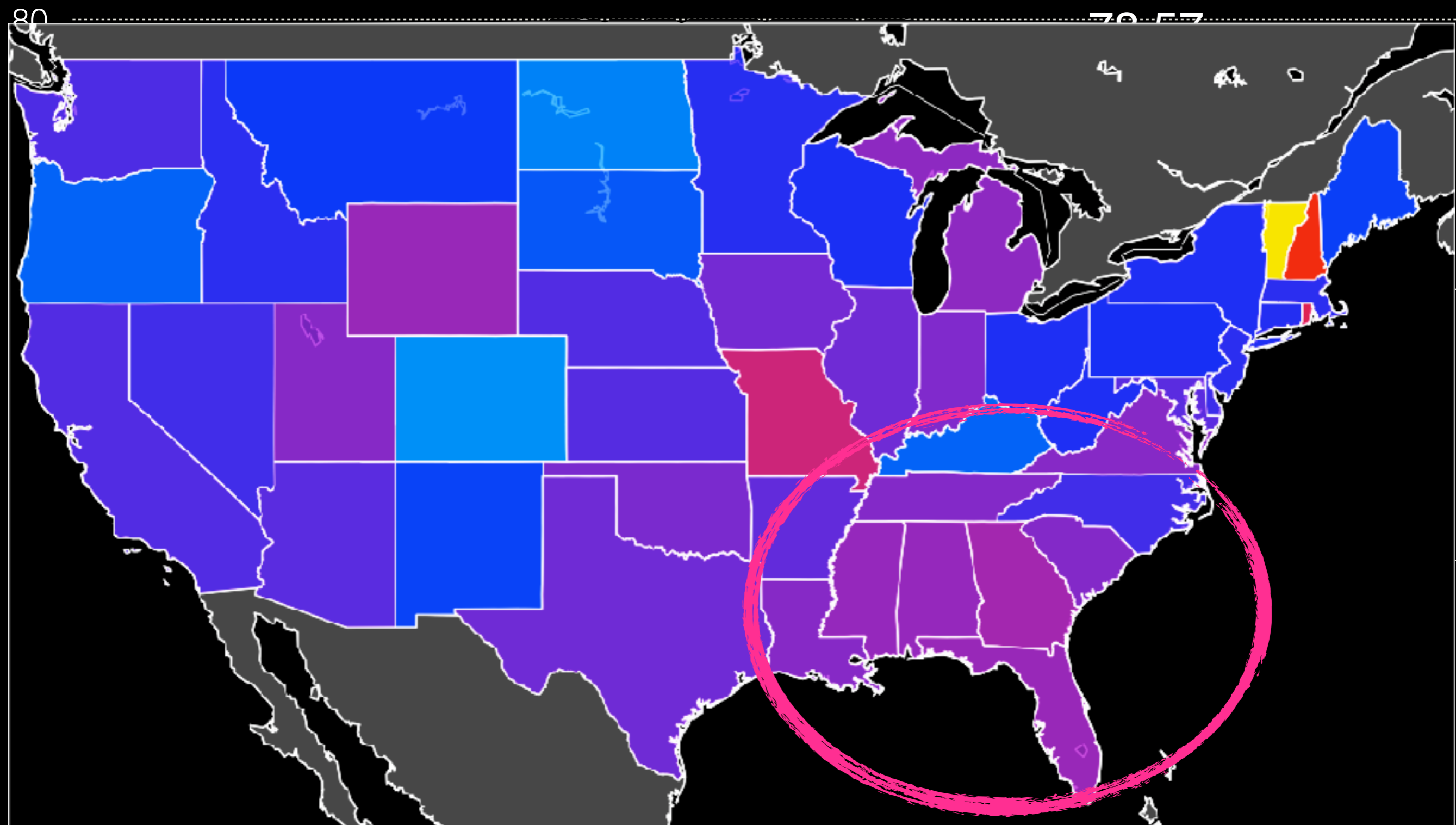
It's shite being Scottish in a smart speaker world  
70,140 views

1.7K 118 SHARE SAVE ...



# Biased Tools cause Exclusion

F1



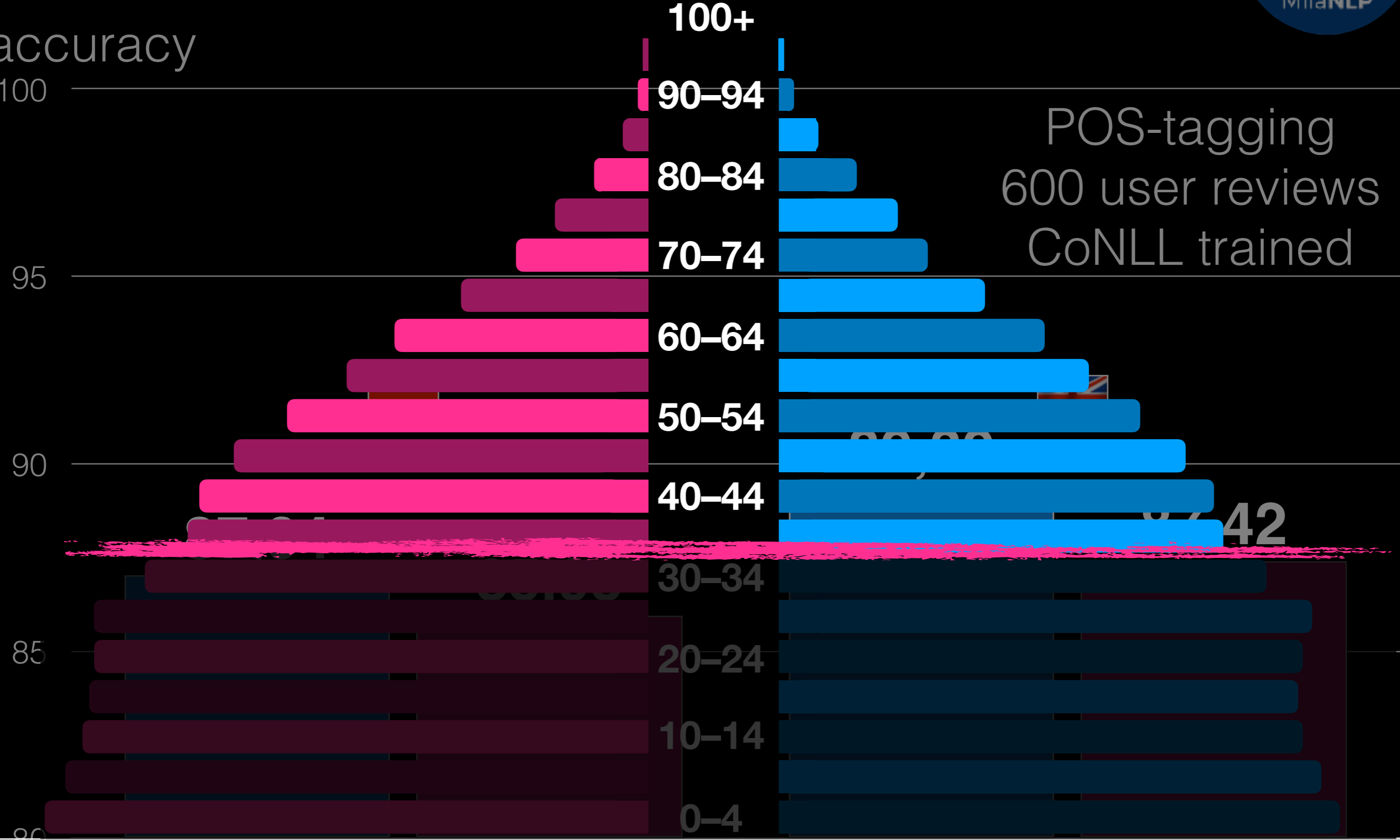
80  
70  
60  
50  
40  
30  
20  
10

Bocconi avg



# Biased Tools cause Exclusion

accuracy  
100



POS-tagging  
600 user reviews  
CoNLL trained

O45

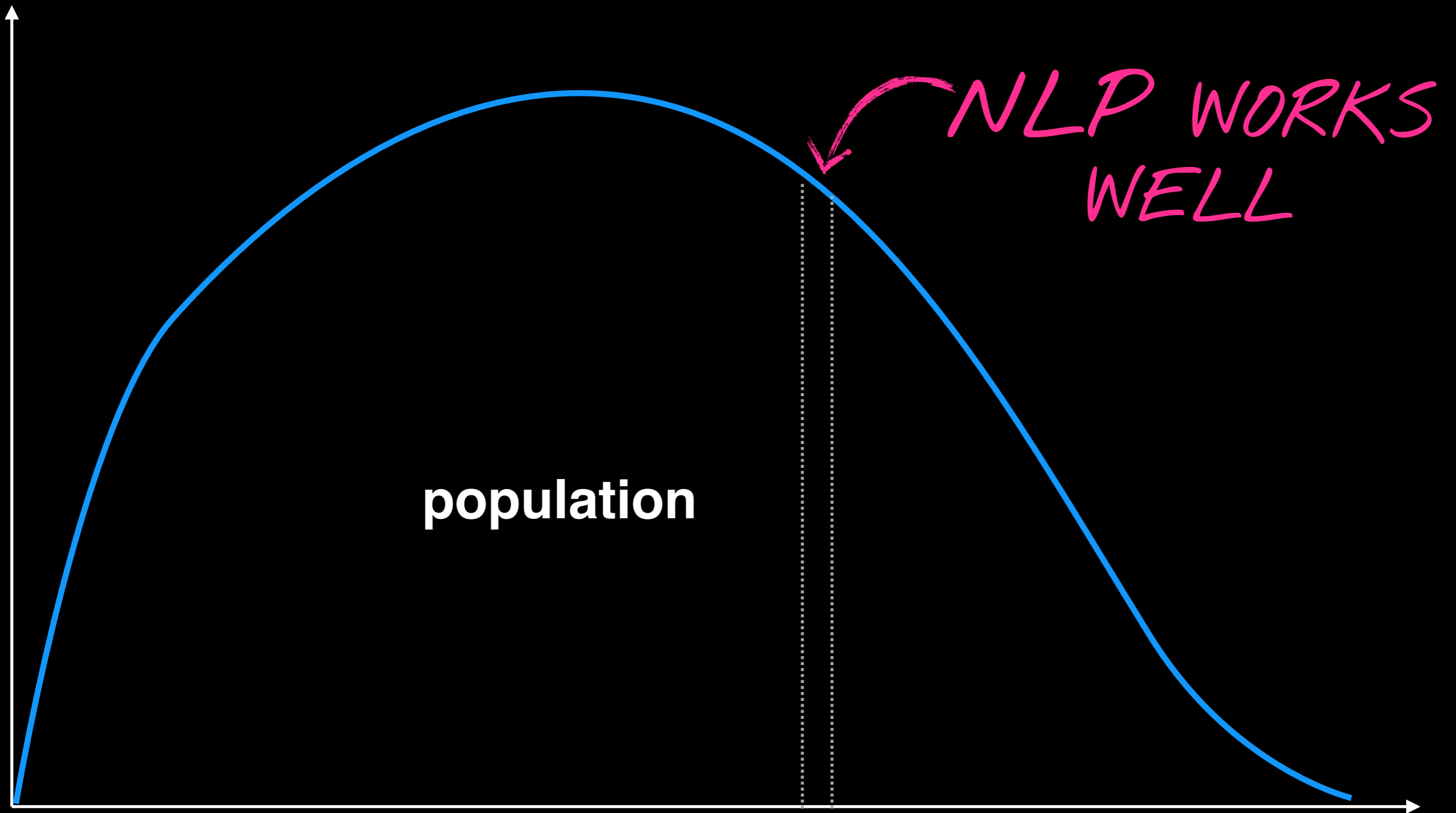
U35

O45

U35



# The Consequences





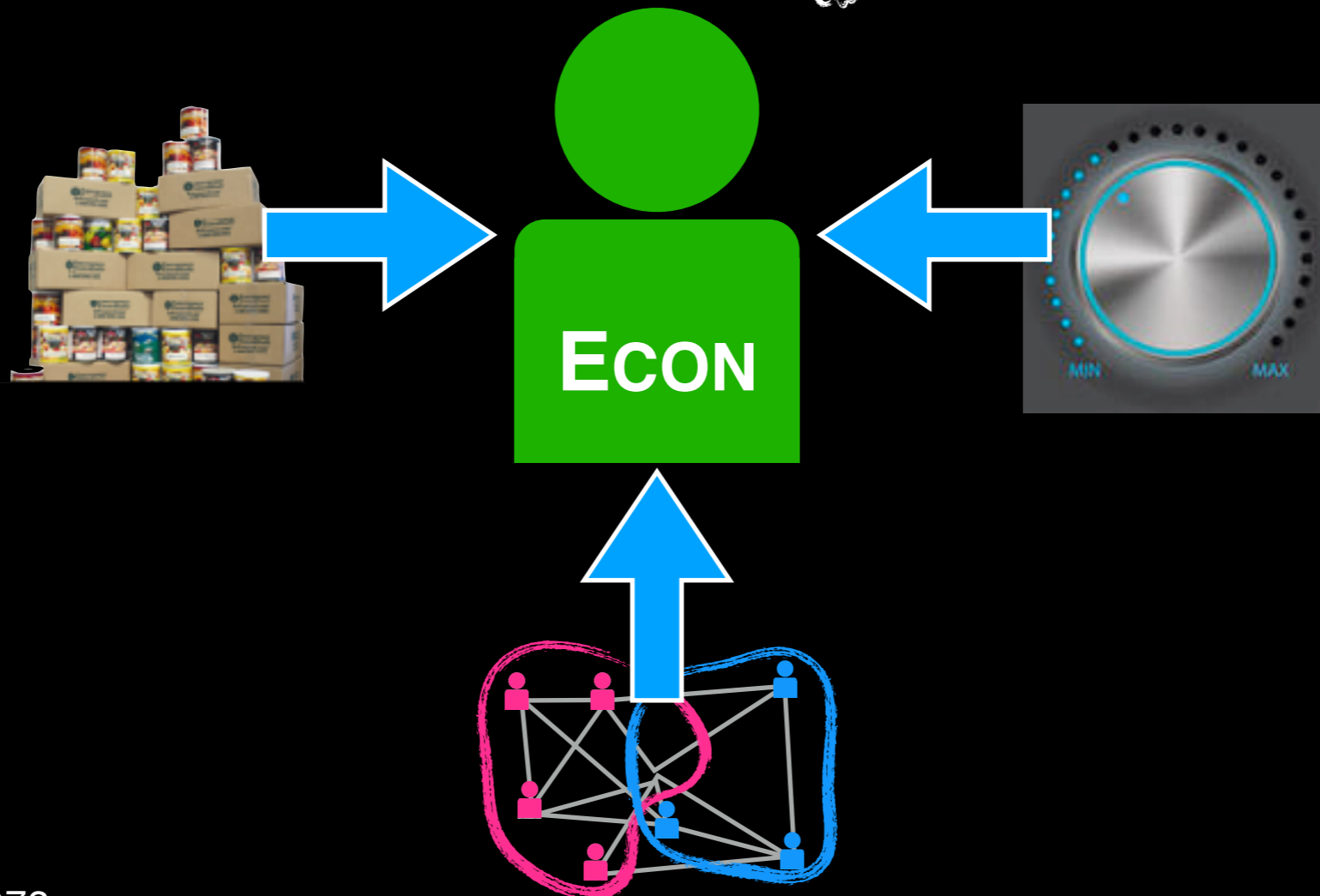
**Where are We?**

# A Limiting View



Daniel Kahneman Amos Tversky

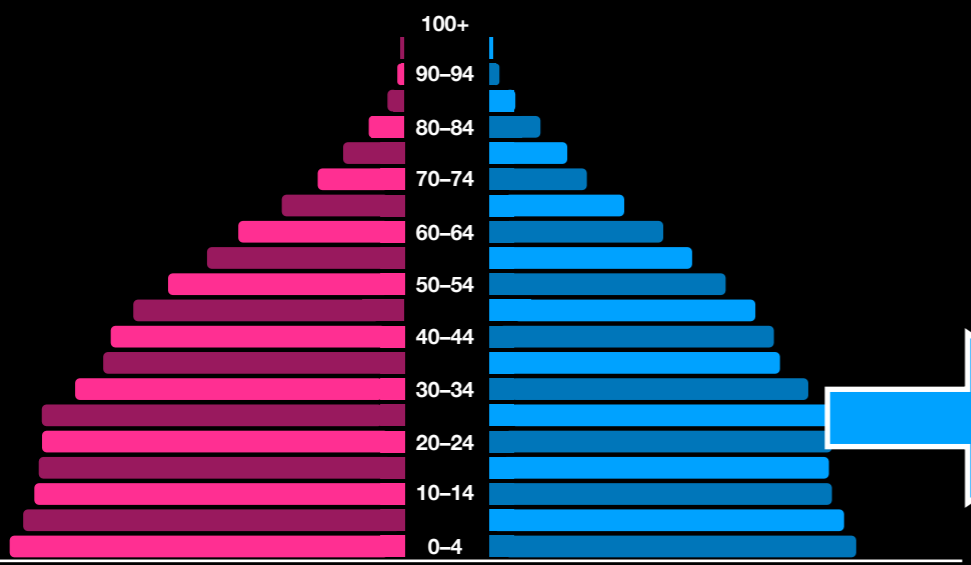
$$y = wX + b$$



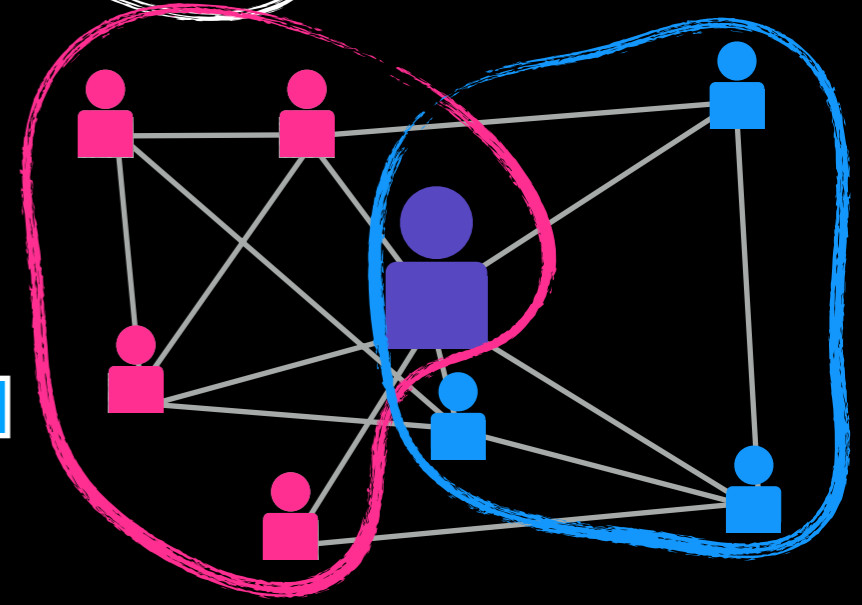
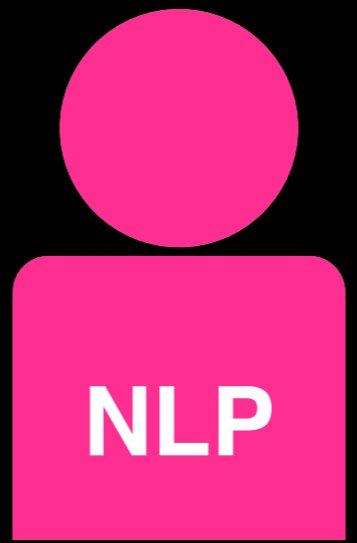
# Language as *SOCIAL* INFORMATION CONSTRUCT



$$y = F(wX + b)$$



*DEMOGRAPHICS*



*SOCIAL*



“[T]he common misconception [is] that language use has primarily to do with words and what they mean. It doesn’t. It has primarily to do with people and what *they* mean.”

*Clark and Schober (1992)*



# Social Factors



communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

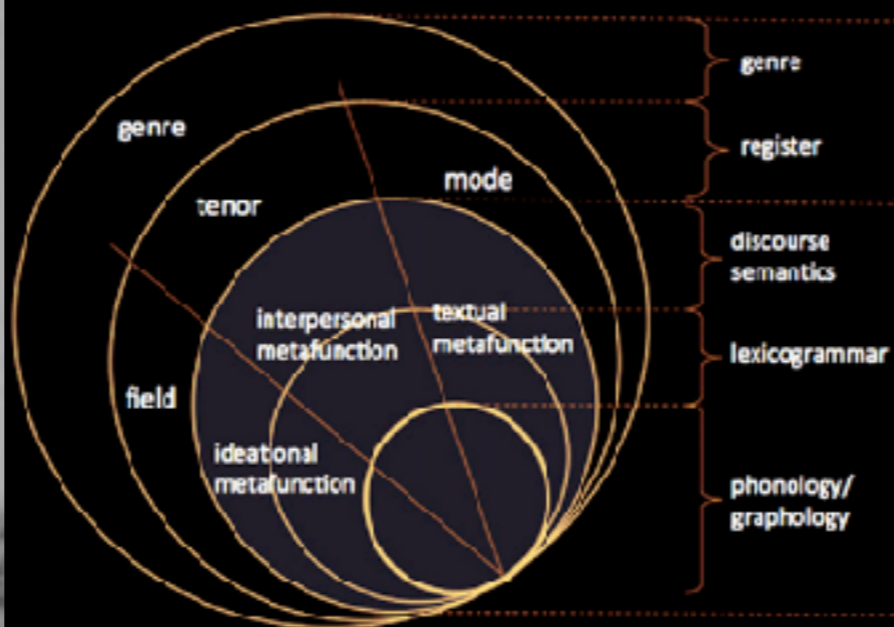
# (Linguistic) Background



"Language games"

Ludwig Wittgenstein

## Systemic Functional Linguistics



"Cooperative Principle"

Herbert Paul Grice

# Speaker



communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

Consistent traits  
signaling demographic  
makeup. Improve NLP  
performance

Applications:

Text generation

Text classification

POS tagging

Conversational agents

35–  
40

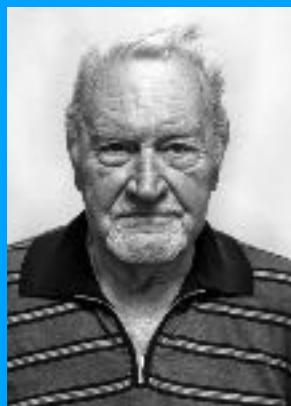
male

academic

NLP

in Milan

# Implicit Bias



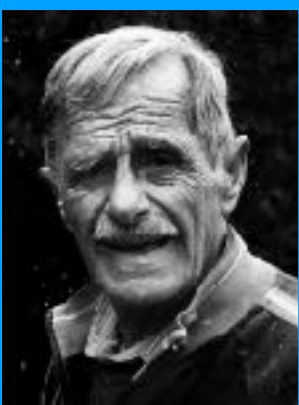
Example 1



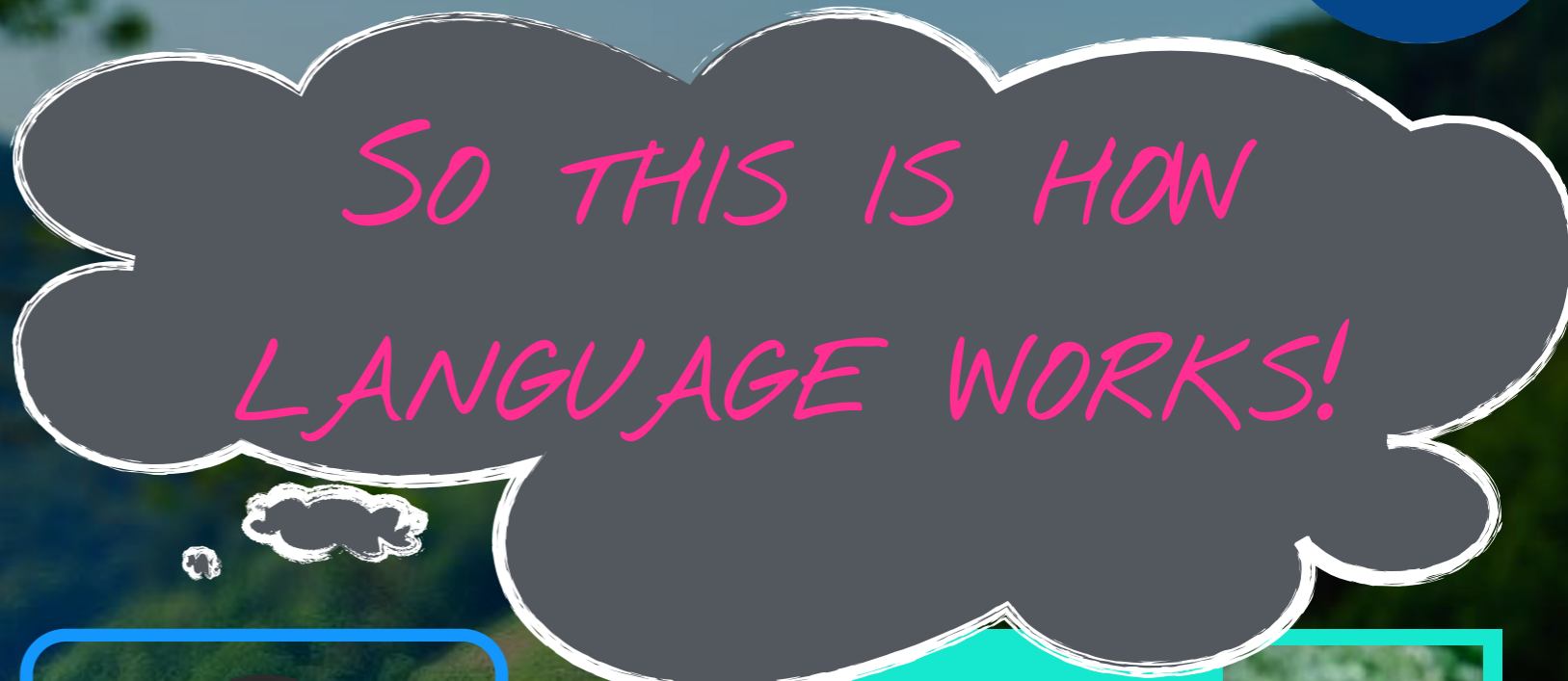
Example 2



Example 3



Example N



Hello,  
computer



You sound different...

# Receiver



communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

Intended audience and their expectations.

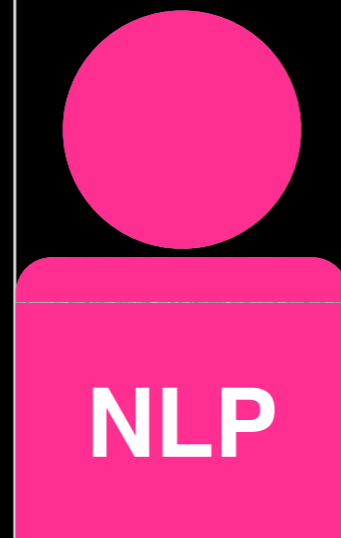
Applications:

Text generation

Conversational agents

Spellchecking

Hey, can't make tonite,  
sorry!



Dear Madam President,  
I regret to inform you  
that I will not be able  
to participate.

communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

# Applications

**Where  
does the sun go  
at night?**

**Here is "Of the  
movement of celestial bodies"  
by Newton.**

# Social Relation

communicative goals	
culture & ideology	
social norms	
context	
social relation	
speaker	receiver

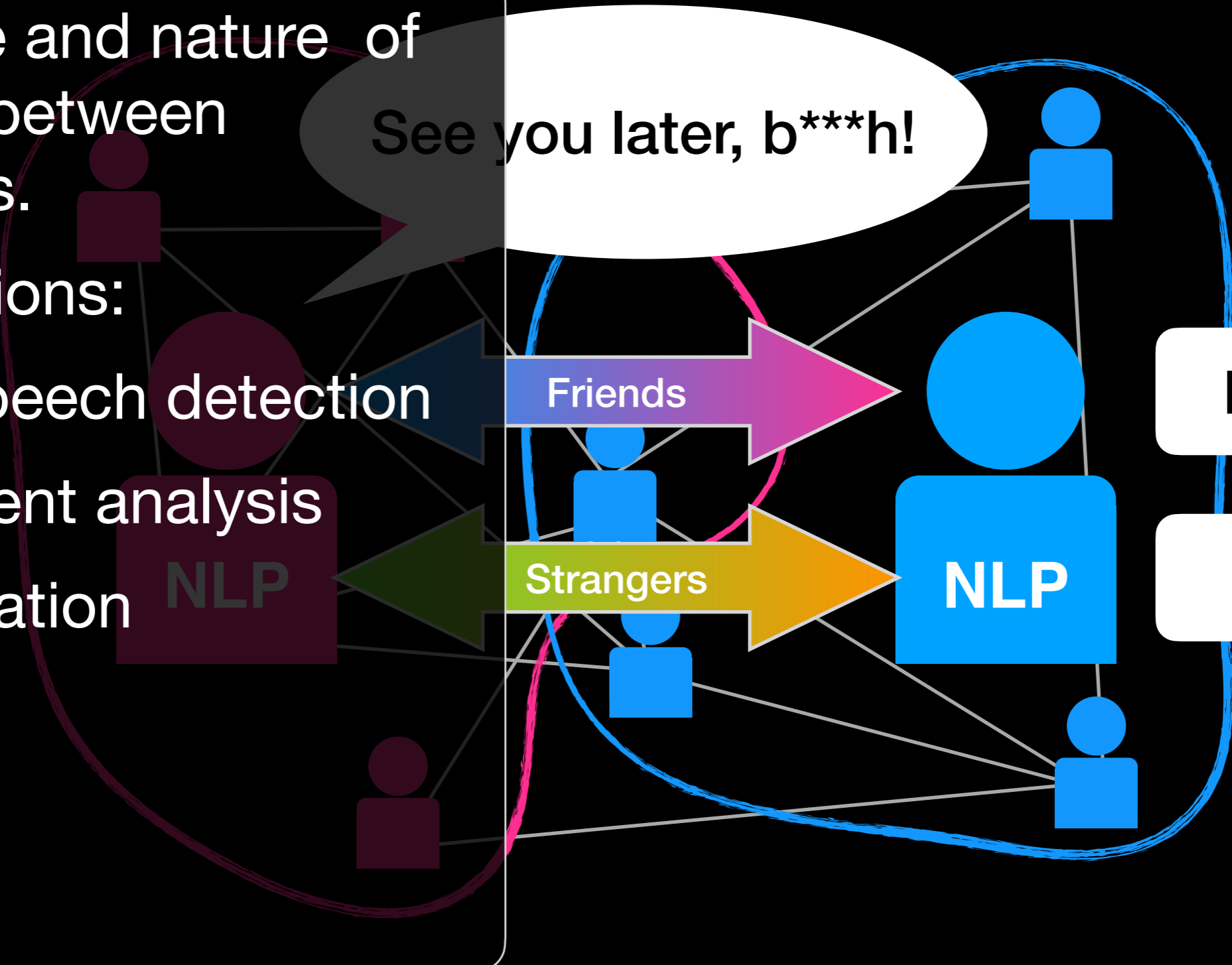
Distance and nature of relation between speakers.

Applications:

Hate speech detection

Sentiment analysis

Geolocation



communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

# (Social) Context



Non-textual factors:  
domain, language,  
situation, topic, etc.

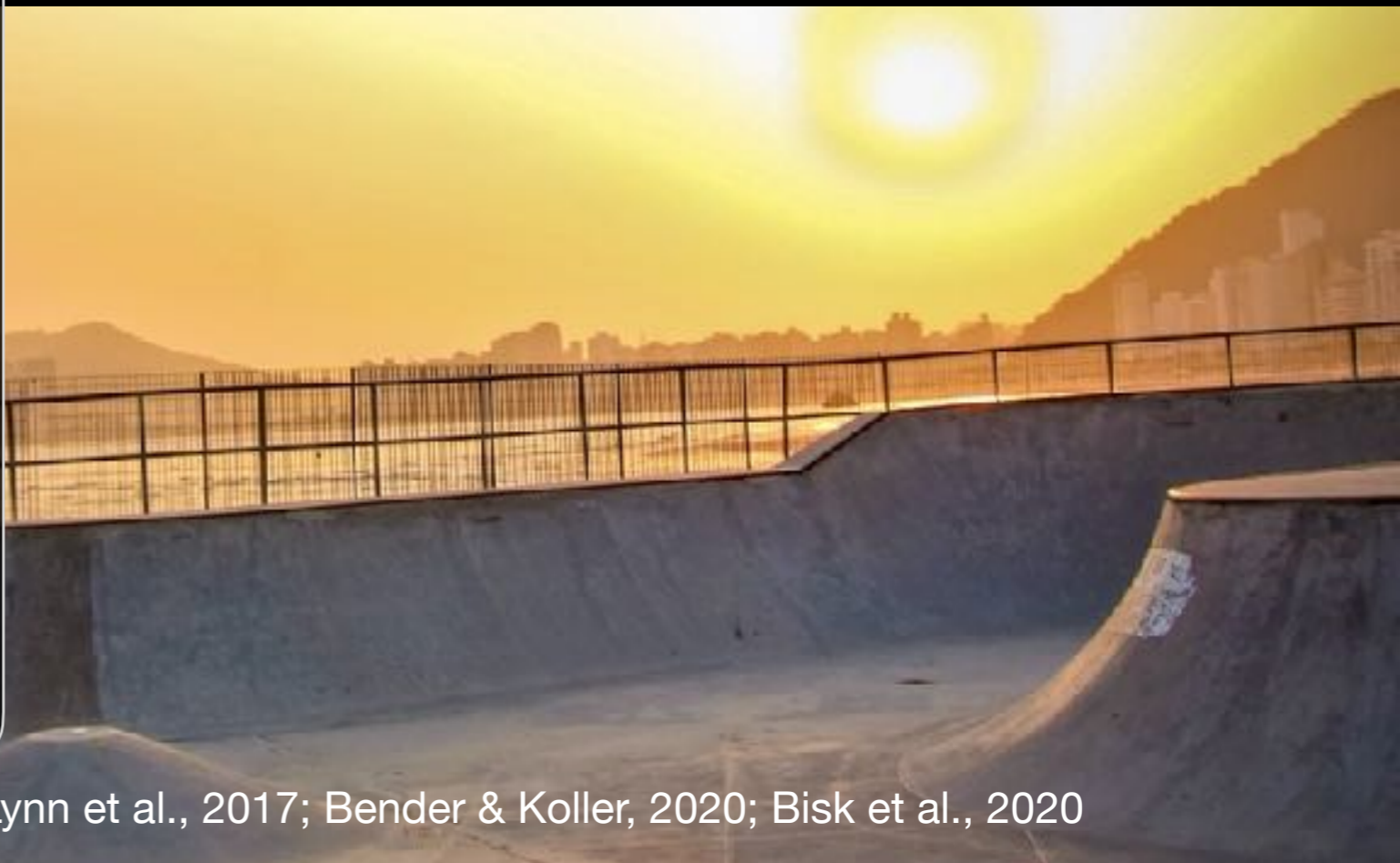
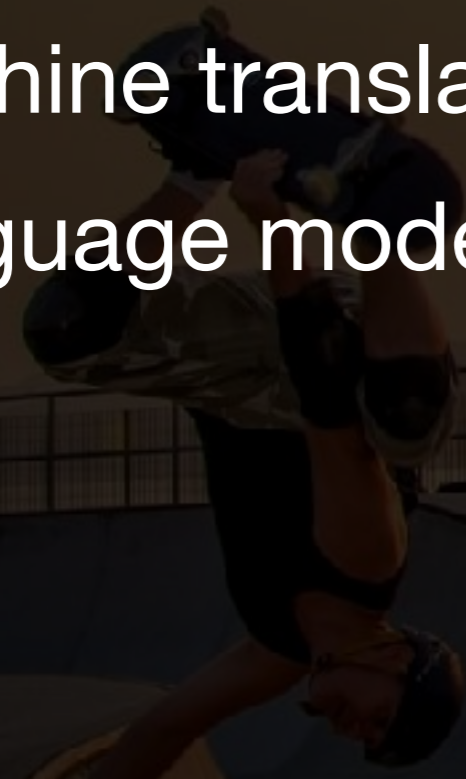
Applications:

Text generation

Machine translation

Language models

That was a **sick performance!**





# Language Invariant Properties



$$f(a) = f(T(a))$$

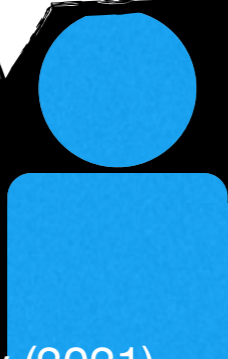
LINGUISTIC PROPERTY THAT STAYS THE SAME  
EVEN WHEN THE TEXT IS TRANSFORMED.

That was a sick  
performance!

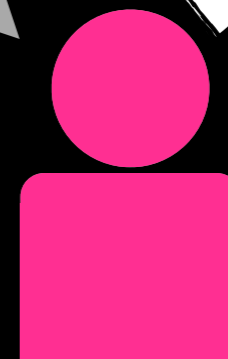


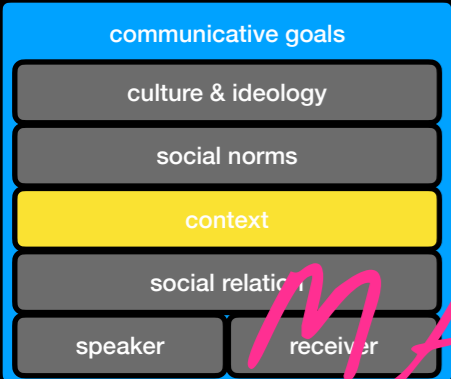
Translate

É stata  
una prestazione  
disdicevole! X



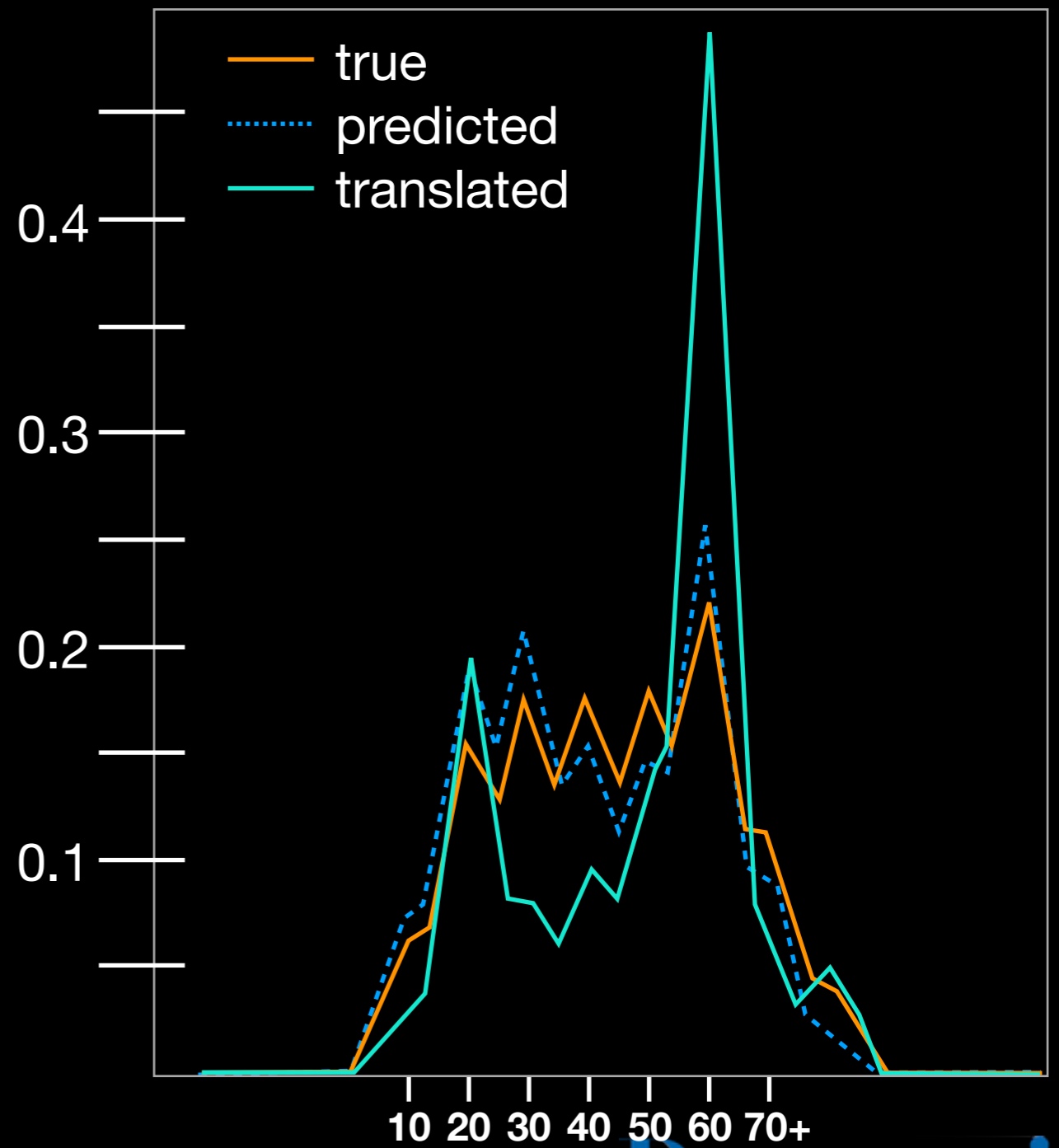
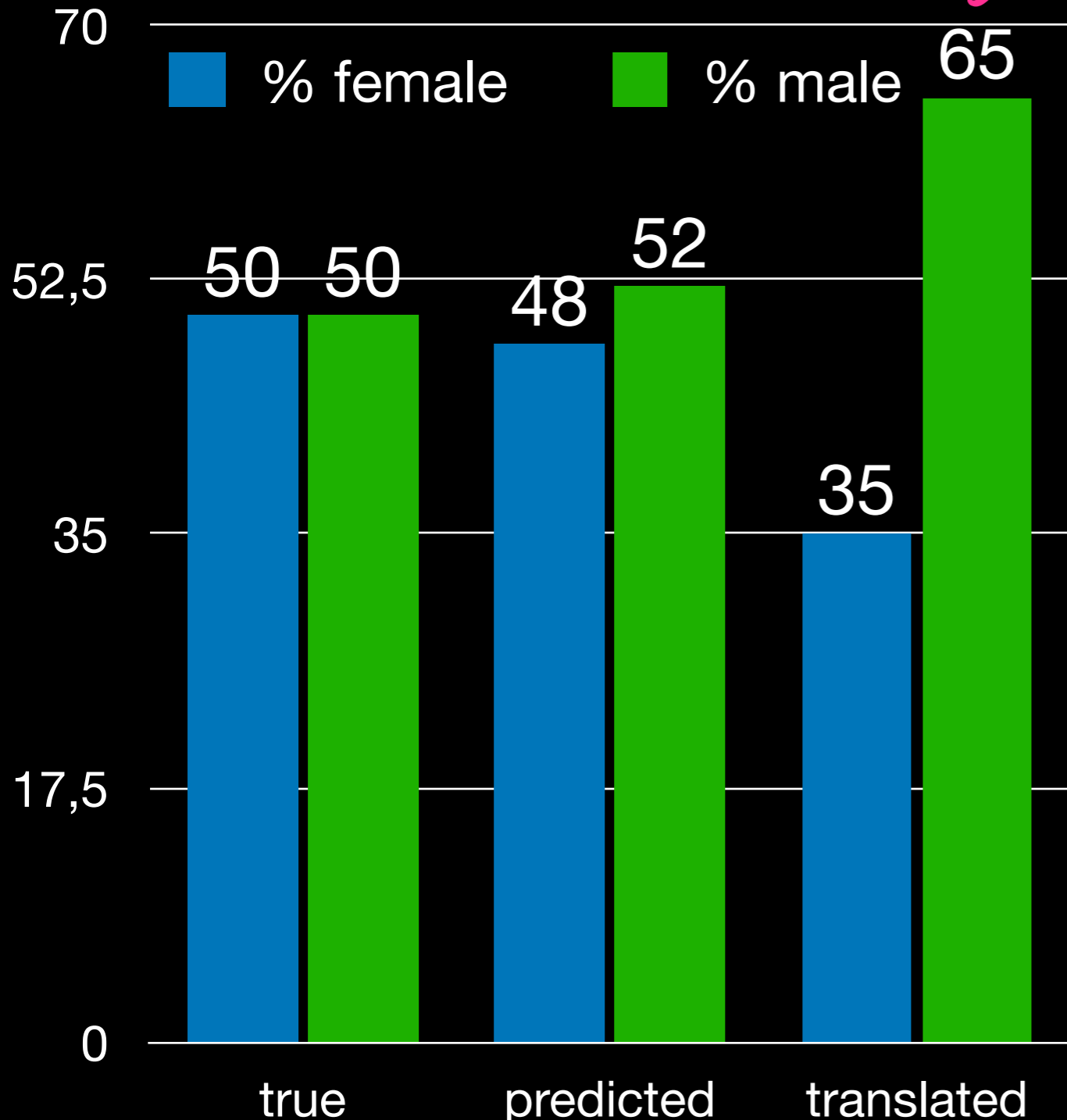
É stata una  
prestazione super! ✓





# Applications

*MACHINE TRANSLATION MAKES YOU SOUND OLDER AND MORE MALE.*



communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

# Social Norms



Acceptable conduct,  
shared understanding.

Applications:

Content moderation

Conversational agents

Annotation models

I love you!

I know...

communicative goals

culture & ideology

social norms

context

social relation

speaker

receiver

# Culture and Ideology



Customs, ideology, and cultural identity.

Applications:

Text generation

Conversational agents

Machine translation

I don't know!

Where's the pharmacy?

Take the first left, walk 5min, go right at the intersection, and keep on going til you see a large tree, then...

*INCORRECT, BUT POLITE*

# Communicative Goals



Metafunction of ideational and interpersonal goals.

Applications:

Text generation

Conversational agents

Humor

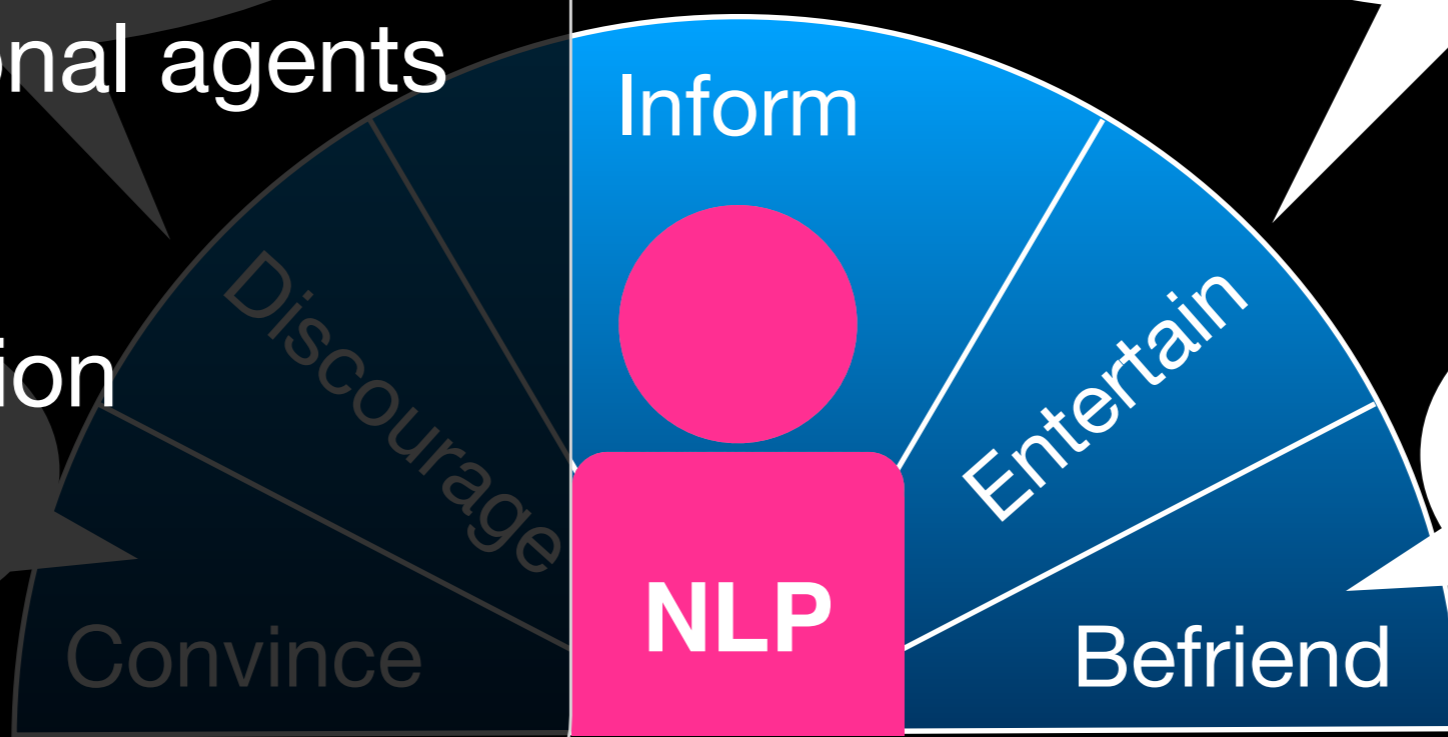
Argumentation

get it done in no time!

You and I will write the report.

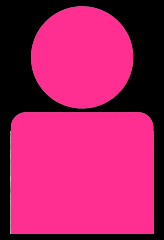
Guess who gets to write the report?

I'll help you with the report!





**Aside:**  
***Do Pay Attention to the  
Man behind the Curtain!***



# Annotator Bias



Whatever,  
it's **X**

No! It's a  
**NOUN!**



*HAS NO CLUE...*

PRON VERB ADP

**X**

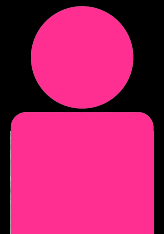
NOUN

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

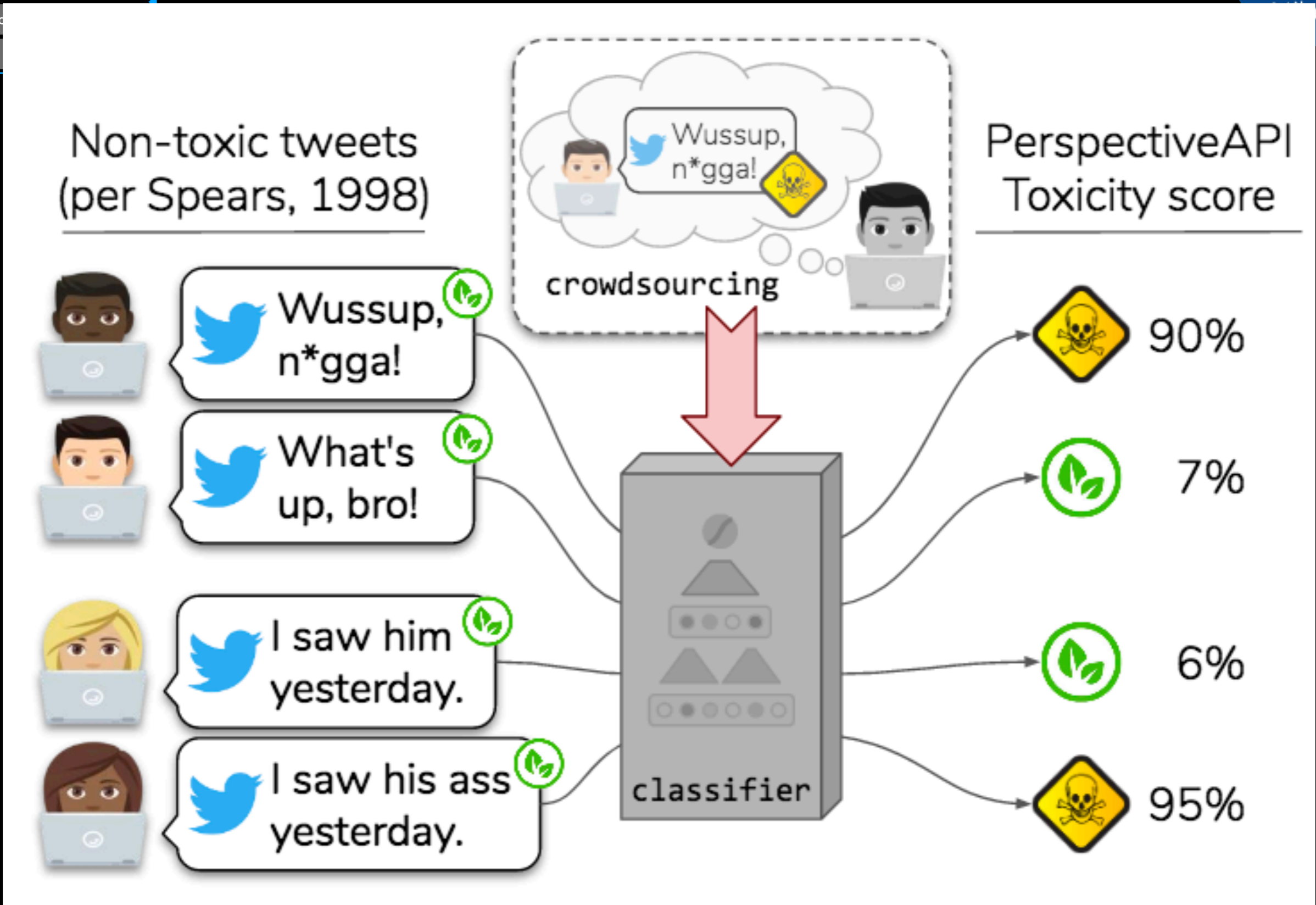
**NOUN**

NOUN

it is on social media



# Annotator Bias...





European Research Council  
Established by the European Commission



# What to Do Next?

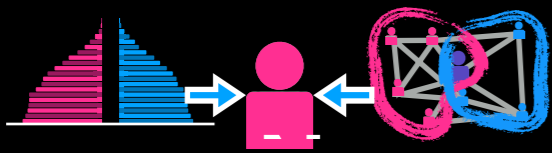


**INTEGRATOR**  
European Research Council

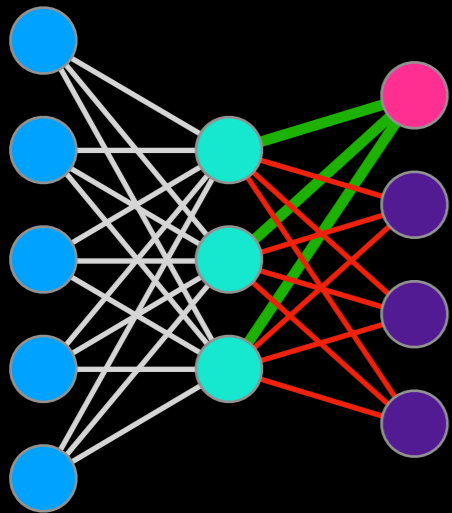
# INTEGRATOR Objectives



1. provide data sets and metrics



2. provide theoretical foundations



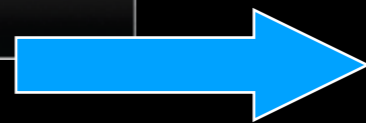
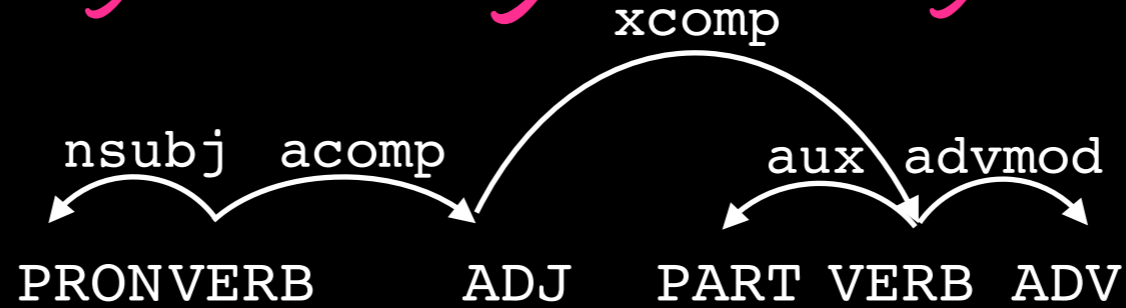
3. algorithmically integrate demographic knowledge



# Better Data

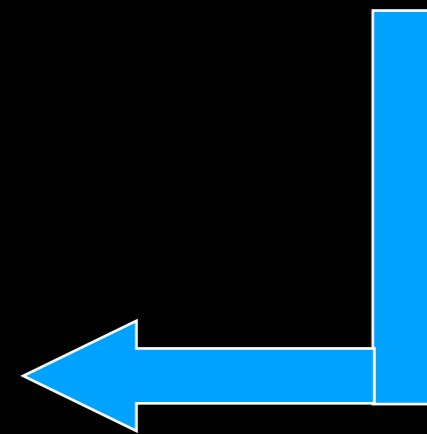
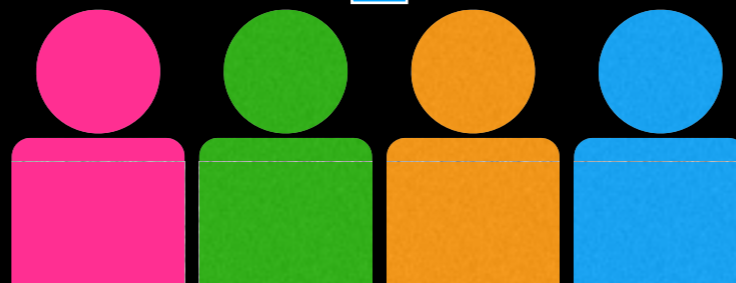


## SYNTACTIC ANNOTATIONS



ANNOTATE,  
INFER  
DEMOGRAPHICS

DOCUMENT  
DEMOGRAPHICS!

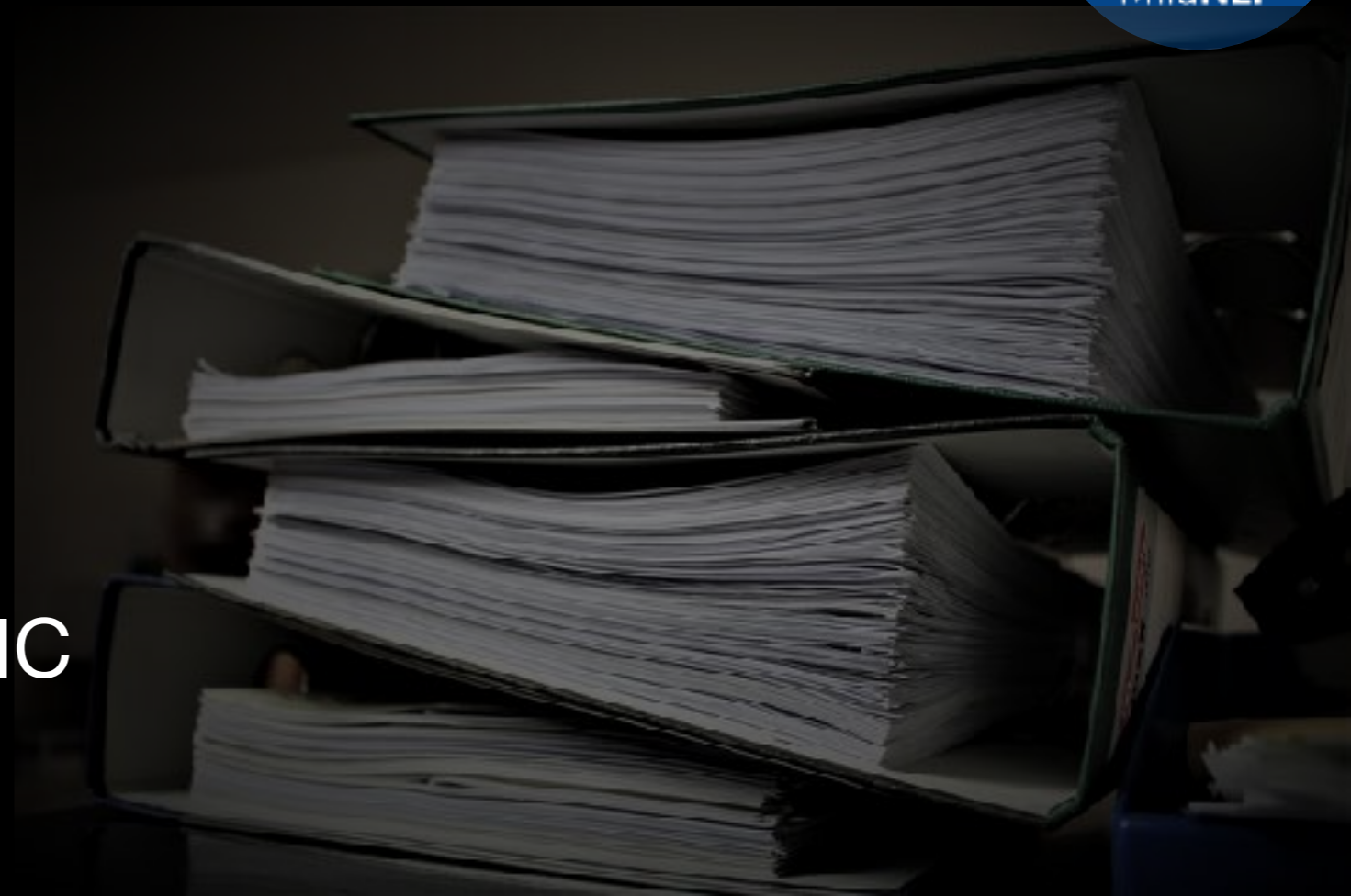




# Data Statements



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



*RECONSTRUCT  
COLLECTION*



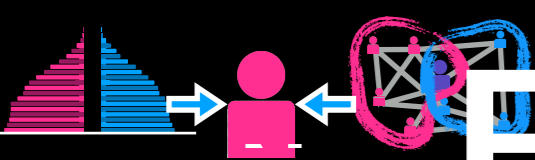
# The Bender Rule



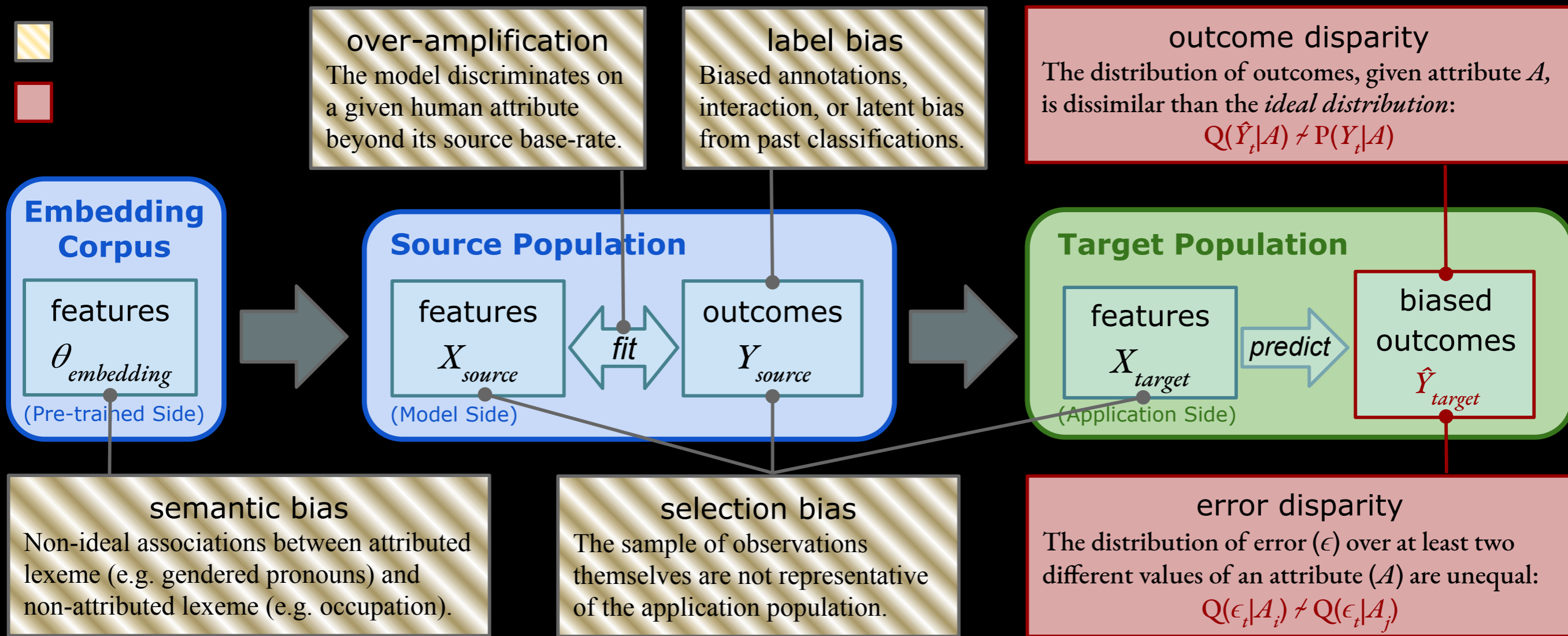
***"Do state the name of the language that is being studied, even if it's English."***

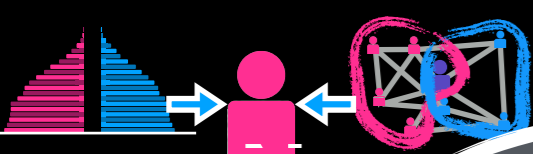
***– Emily Bender***



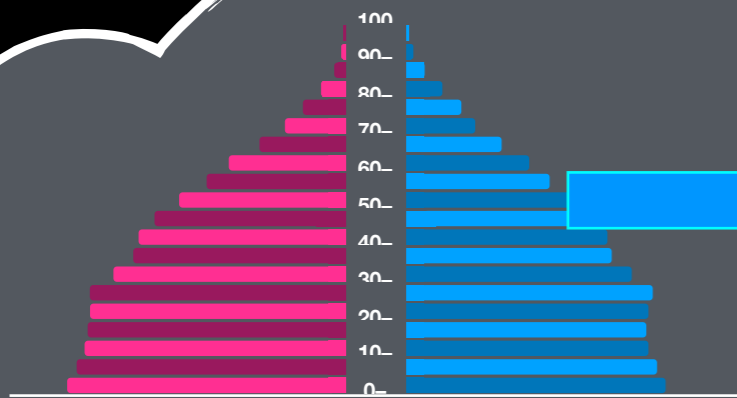


# Bias Taxonomies

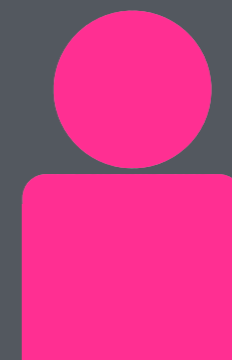




# Sources of Bias



*DATA*



*ANNOTATION*



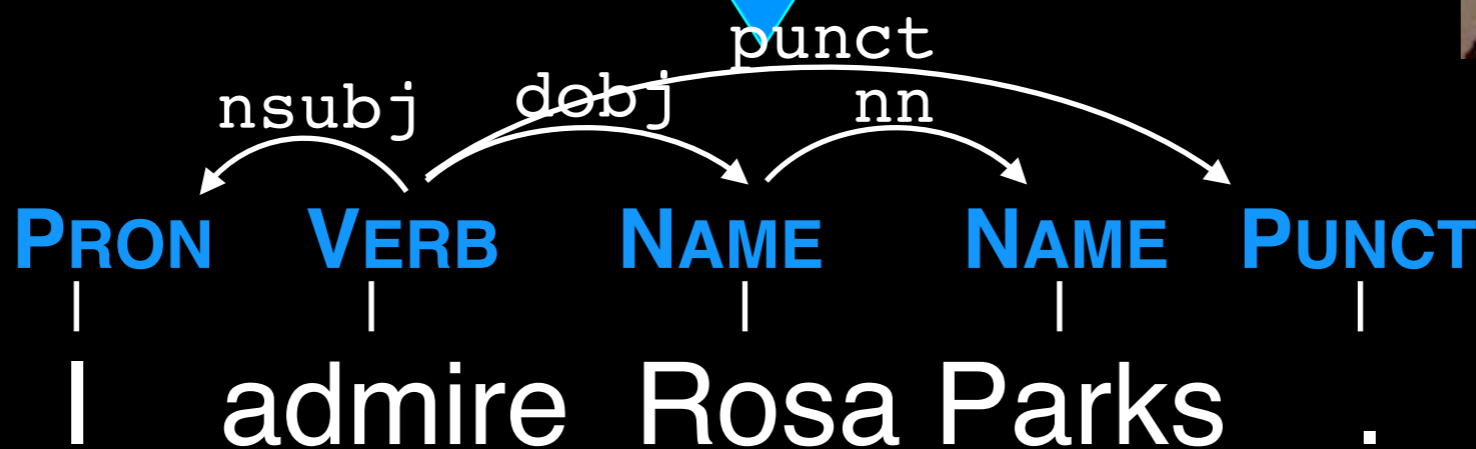
*REPRESENTATIONS*



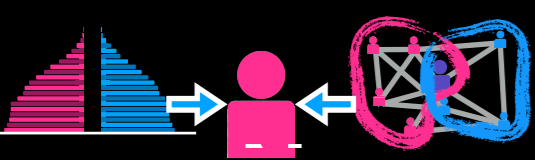
*MODELS*



*DESIGN*







# Better Metrics



La **donna** è una gran tr\*\*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*

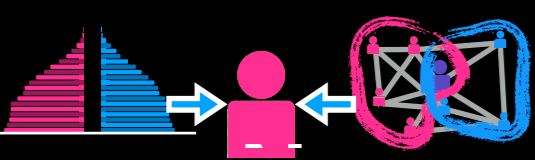
La **mujer** es la mas bonita

*The wife is the most attractive*

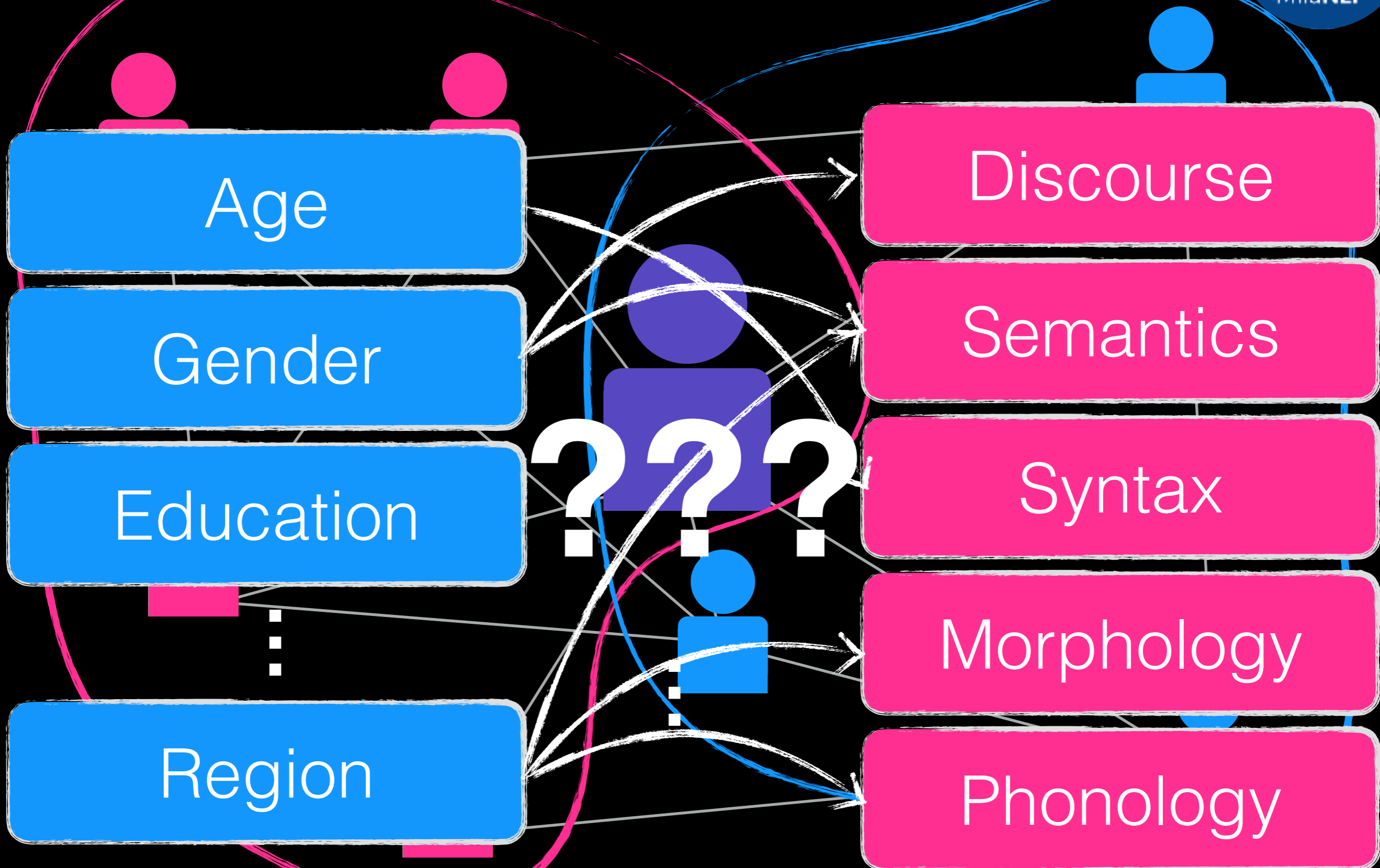
**4.5%** of all completions contain a hurtful word.

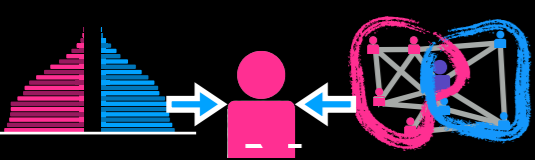
if the target inflection is **female**, **10%** refer to **sexual promiscuity**

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

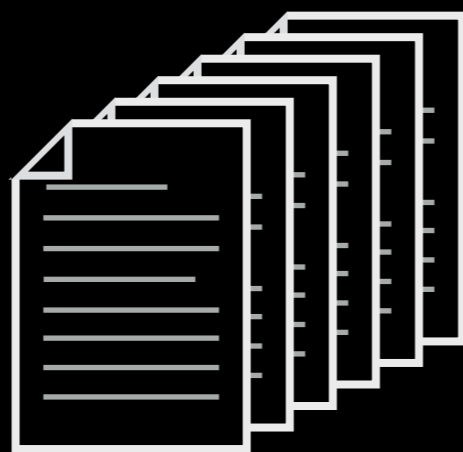


# What to Model?





# Broad Domain...



WSJ

*DOMAIN*



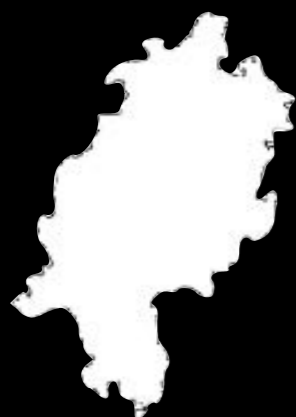
Barbara Plank

O45

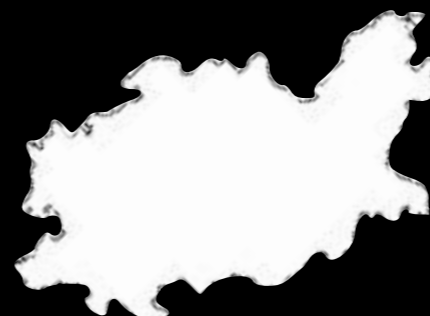
*AGE*



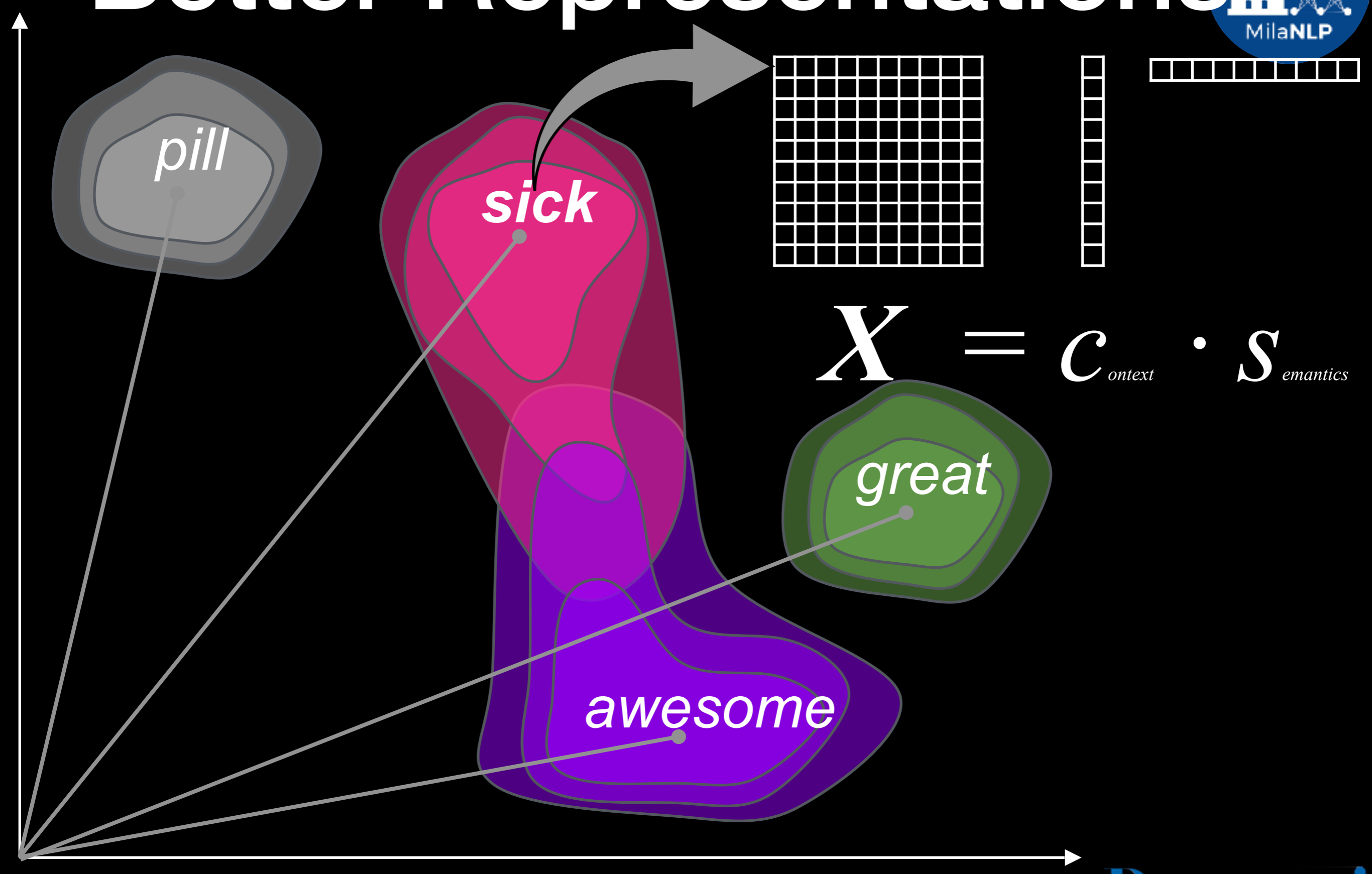
U25

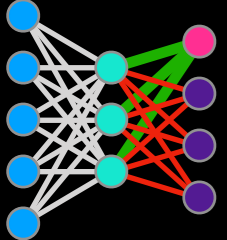


*DIALECT*



⋮

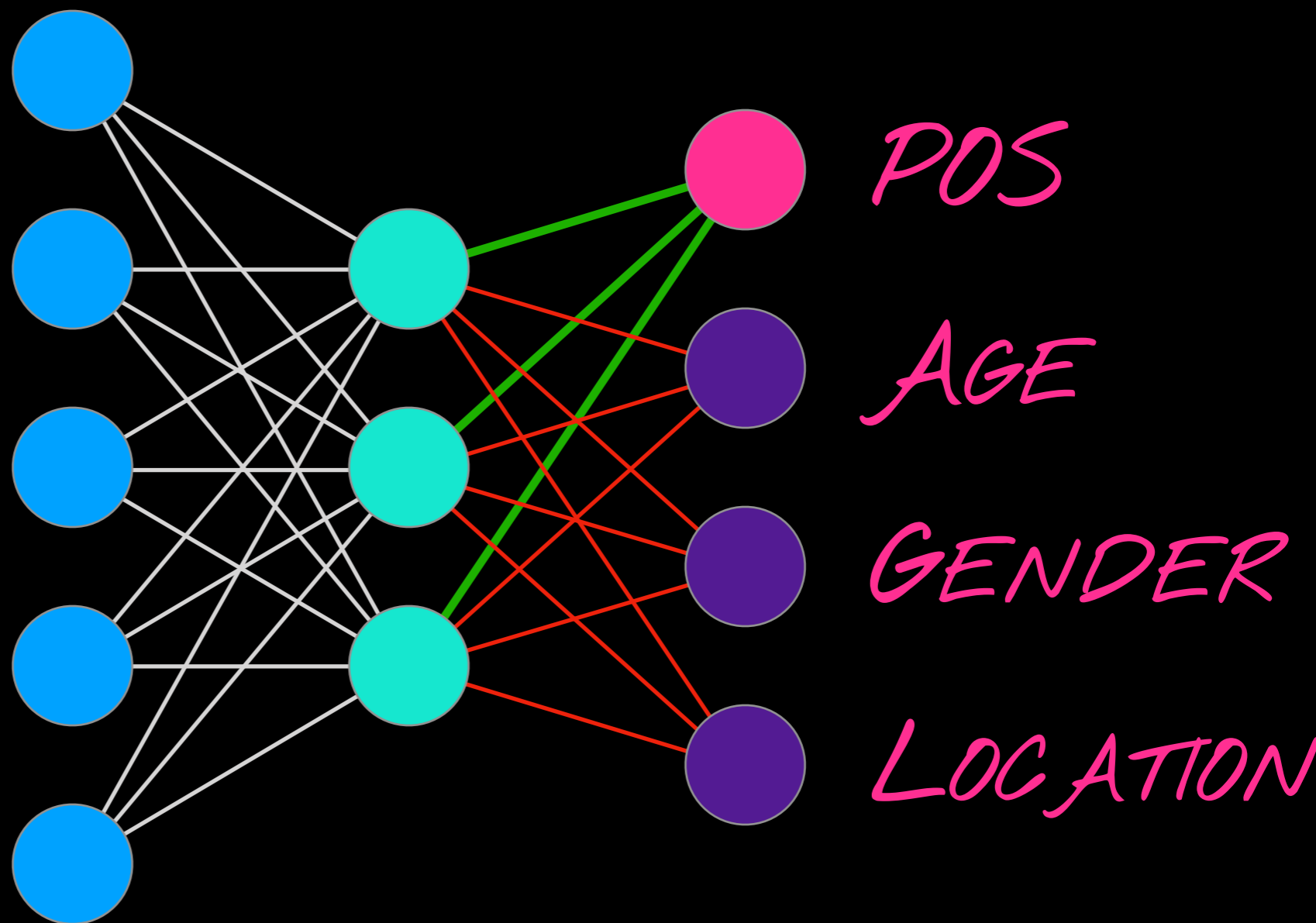




# Better Models



- representation learning
- domain adaptation
- multitask learning
- reinforcement learning



# 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

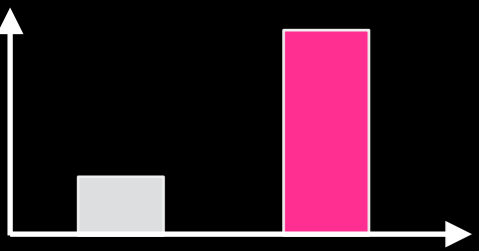


# Wrapping Up

# Social Factor Goals



- develop novel methods to include social factors into NLP



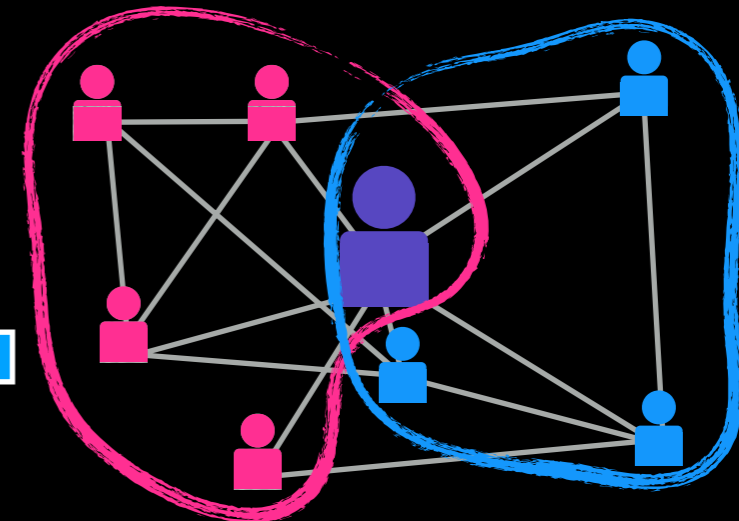
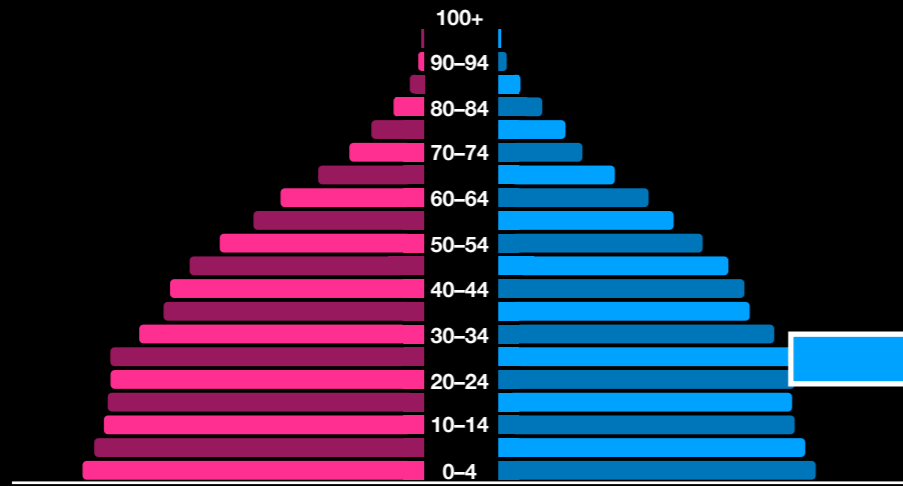
- improve existing performance



- enable novel personalized applications



# Take-Home Points



*DEMOGRAPHIC*

*SOCIAL*

- People use language for information *and* social purposes
- NLP has modeled only information
- We are missing out on performance, personalization, fairness
- Let's model social aspects of language!



[www.dirkhovy.com/portfolio/papers](http://www.dirkhovy.com/portfolio/papers)  
[milan1proc.github.io/](https://github.com/milan1proc)

# Thank you!

*IF YOU ARE INTERESTED IN WORKING  
WITH ME: I HAVE UPCOMING POSTDOC  
POSITIONS IN THESE AREAS*

 @dirk\_hovy

[www.dirkhovy.com](http://www.dirkhovy.com)



[www.dirkhovy.com/portfolio/papers](http://www.dirkhovy.com/portfolio/papers)  
[milan1proc.github.io/](https://github.com/milan1proc)

# Questions?

*IF YOU ARE INTERESTED IN WORKING  
WITH ME: I HAVE UPCOMING POSTDOC  
POSITIONS IN THESE AREAS*

 @dirk\_hovy

[www.dirkhovy.com](http://www.dirkhovy.com)