

# Enhancing Hate Speech Annotations with Background Semantics

**Paula Reyero Lobo**

Knowledge Media Institute  
The Open University, UK

23 October 2024

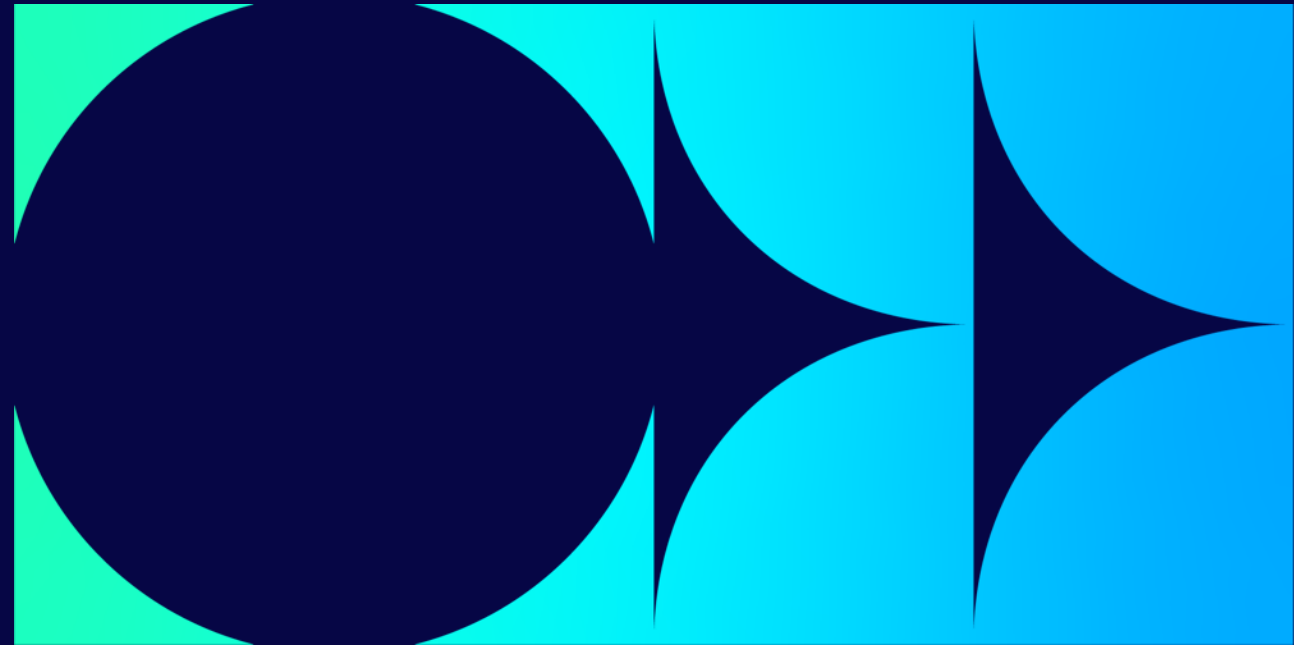
Advised by:

Enrico Daga, Harith Alani, Miriam Fernandez



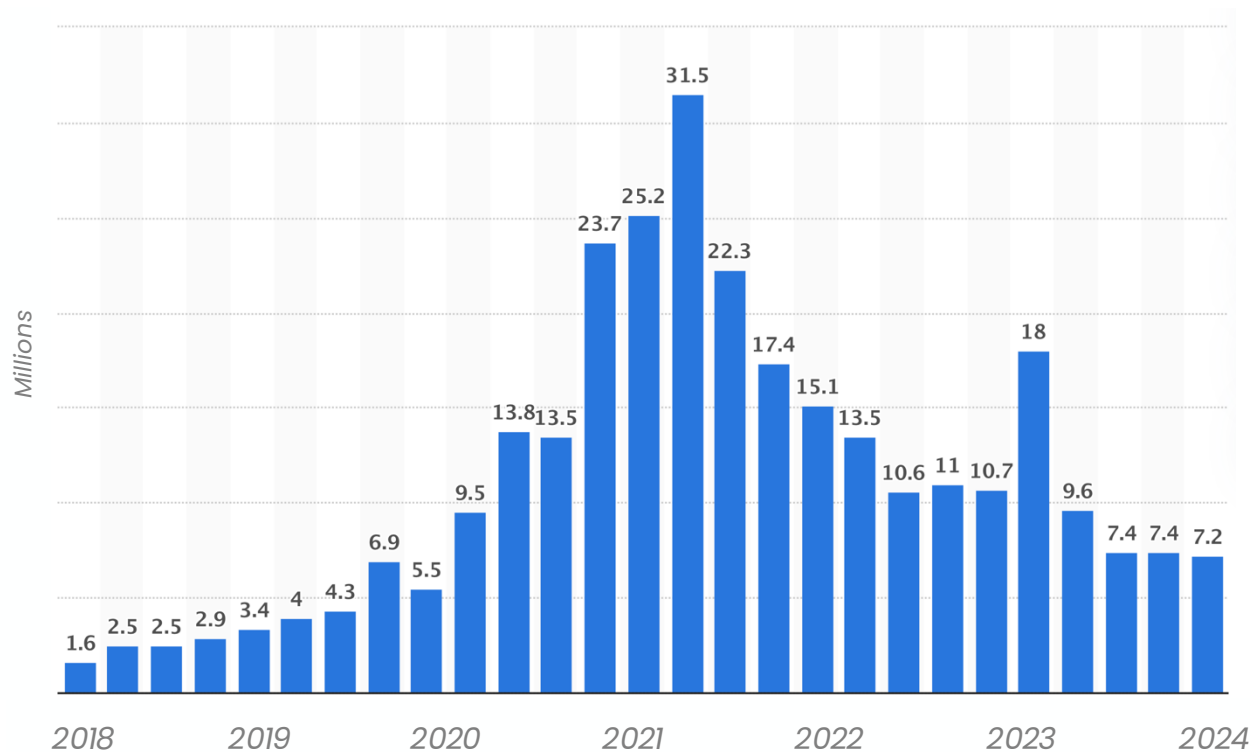
Funded by  
the European Union

**Content warning.** This research aims to tackle hate speech and contains examples of triggering and harmful language.



# Introduction

Content in millions removed on Facebook



Statista research Department, 12 September 2024

# Introduction

Does this comment contain **hate speech**,

defined as "*bias-motivated, hostile and malicious language targeted at a person/group because of their actual or perceived innate characteristics*"?



The kebabs are a bunch of homosexual rapist deviants.



Suck farts out of obese sheboons



Does the message mention or is about...

**gender**  
**sexuality**  
**race...?**

Human annotations are crucial for training hate speech detection models

# Introduction

Does this comment contain **hate speech**,

defined as "*bias-motivated, hostile and malicious language targeted at a person/group because of their actual or perceived innate characteristics*"?

Does the message mention or is about...

**gender**  
**sexuality**  
**race...?**



The kebabs are a bunch of homosexual rapist deviants.



**Yes**

**No**

Label: 0.5



Suck farts out of obese sheboons



**Woman**

Label: 1.0



The disagreement between annotators challenges the use of majority votes and the accuracy of AI models

## Related work

Does this comment contain **hate speech**,

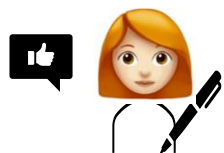
defined as "*bias-motivated, hostile and malicious language targeted at a person/group because of their actual or perceived innate characteristics*"?

Does the message mention or is about...


**gender**  
**sexuality**  
**race...?**

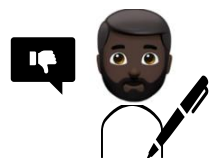


The kebabs are a bunch of homosexual rapist deviants.



No

Label: 0.0 

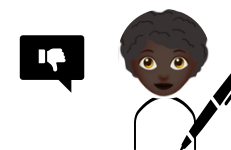


Yes


Label: 1.0 



Suck farts out of obese sheboons



Woman

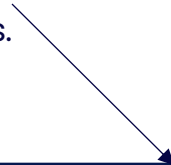
Label: 1.0 



# Problem statement



The **kebabs** are a bunch of **homosexual rapist** deviants.

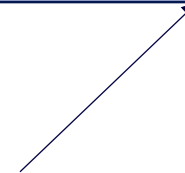


Ethnic slur to refer to a **Muslim**



Suck farts out of obese **sheboons**

Slang, derogatory, offensive term targeting **Black women**



The use of slang, non-standard, and domain specific language makes groups targeted by hate speech especially difficult to understand and perceive.

# Our approach

To compensate for the lack of familiarity with domain language, we introduce semantics to enhance comprehension in hate speech annotation tasks.



Data Repository

A post:

Post:

"ana navarro is a spic cunt that needs to shut her commie cock holster, these pieces of shit declare war on us every time they open their ugly commie mouths how long until enough is enough"

Part 1: Identify if there are any references to gender and/or sexuality.

A post with **background semantics**: highlighted entities with definitions

Post:

ana navarro is a [spic cunt](#) that needs to shut [her](#) commie [cock](#) holster, these pieces of [shit](#) declare war on us every time they open their ugly commie mouths how long until enough is enough

Highlighted definitions:

'spic'

'A highly offensive racial slur referring to any Latino/Hispanic person.'

Part 1: Identify if there are any references to gender and/or sexuality.

Does the message mention or is about gender?

Select any or all that apply.

- ☐ Men ☐ Women ☐ Non-binary ☐ Other gender
- ☐ It specifically mentions or is about transgender.

# Our approach

Phase 1

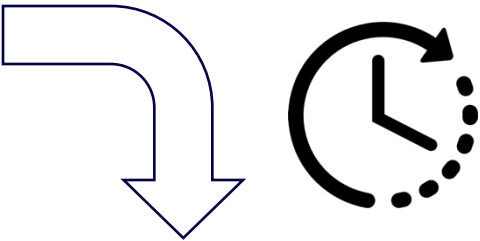
IntroductionPrevious ExperienceDemographic InformationAnnotation GuidelineContent Warning

Previous123456789101112131415SubmitNext

Post:

*"ana navarro is a spic cunt that needs to shut her commie cock holster, these pieces of shit declare war on us every time they open their ugly commie mouths how long until enough is enough"*

One week later



Phase 2

IntroductionAnnotation GuidelineContent Warning

Previous123456789101112131415SubmitNext

Post:

*ana navarro is a spic cunt that needs to shut her commie cock holster, these pieces of shit declare war on us every time they open their ugly commie mouths how long until enough is enough*



Data Repository



# Our approach

Phase 1

Introduction Previous Experience Demographic Information Annotation Guideline Content Warning

Previous 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 Submit Next

Post:

"ana navarro is a spic cunt that needs to shut her commie cock holster, these pieces of shit declare war on us every time they open their mouths, these pieces of shit are not enough"

M M W S G G

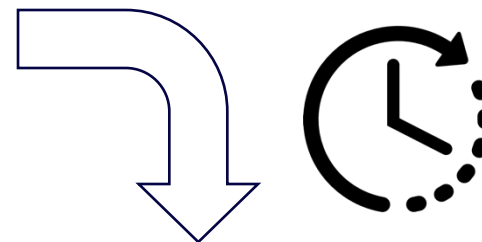
Participants

Each text has annotations from at least three non-cisgender and/or non-heterosexual annotators



Data Repository

One week later



Phase 2

Introduction Annotation Guideline Content Warning

Previous 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 Submit Next

Post:

ana navarro is a [spic cunt](#) that needs to shut [her](#) commie [cock](#) holster, these pieces of [shit](#) declare war on us every time they open their mouths, these pieces of shit are not enough

M M W S G G

Participants

We collect **2,880** annotations from Prolific participants with different genders and sexual orientations, obtaining an average of 6 annotations per text.

# Evaluation

## 1. Agreement

*How supplementing text with semantics affects hate speech annotation agreement?*

## 2. Correlation

*How do semantics impact convergence between annotators from hate speech target and non-target groups?*

## 3. Change

*How does the identification of hate target groups change after introducing semantics?*

# Results

## 1. How supplementing text with semantics affects hate speech annotation agreement?

Krippendorff's Alpha Scores

Gender Labels	Phase 1	Phase 2	$\Delta$
other gender	0.242	0.087	-0.155
non-binary	0.151	0.095	-0.056
gender unclear	0.069	0.035	-0.035
transgender	0.386	0.381	-0.006
gender not-referring	0.187	0.269	<b>0.081</b>
men	0.267	0.396	<b>0.129</b>
women	0.370	0.529	<b>0.159</b>
Sexuality Labels	Phase 1	Phase 2	$\Delta$
asexual	0.229	0.206	-0.023
sexuality unclear	0.086	0.065	-0.021
heterosexual	0.147	0.151	0.004
sexuality not-referring	0.305	0.329	<b>0.024</b>
other sexuality	0.202	0.254	<b>0.053</b>
homosexual	0.597	0.654	<b>0.056</b>
bisexual	0.213	0.295	<b>0.082</b>
General Questions	Phase 1	Phase 2	$\Delta$
Hate speech?	0.318	0.321	0.003
Hate speech targeting?	0.255	0.260	0.005
About sexuality?	0.370	0.409	<b>0.039</b>
About gender?	0.211	0.396	<b>0.186</b>
<b>average</b>	0.256	0.285	0.113

Adding background semantics increased agreement by **11.3% on average** for gender and sexuality groups.

# Results

## 2. How do semantics impact convergence between annotators from hate speech target and non-target groups?

Pearson's Correlation Scores

	Phase 1			
	M	W	S	G
other	<b>0.46</b>	0.38	0.4	nan
non-binary	0.28	0.11	<b>0.44</b>	nan
unclear	0.26	0.07	<b>0.31</b>	nan
transgender	0.4	0.3	<b>0.63</b>	nan
not-referring	0.24	<b>0.28</b>	0.27	nan
men	0.24	<b>0.39</b>	0.3	nan
women	0.44	<b>0.46</b>	0.41	nan

	Phase 2			
	M	W	S	G
other	0.0	<b>0.12</b>	0.06	nan
non-binary	0.08	0.1	<b>0.3</b>	nan
unclear	0.03	-0.04	<b>0.05</b>	nan
transgender	0.5	0.39	<b>0.55</b>	nan
not-referring	<b>0.35</b>	0.26	0.3	nan
men	0.39	0.44	<b>0.56</b>	nan
women	0.6	0.55	<b>0.62</b>	nan

	Phase 1			
	M	W	S	G
asexual	0.31	0.1	nan	<b>0.68</b>
unclear	<b>0.23</b>	0.15	nan	0.17
heterosexual	<b>0.24</b>	0.14	nan	0.23
not-referring	0.33	<b>0.45</b>	nan	0.41
other	<b>0.49</b>	0.19	nan	0.43
homosexual	0.72	<b>0.83</b>	nan	0.66
bisexual	<b>0.49</b>	0.45	nan	0.4

	Phase 2			
	M	W	S	G
asexual	<b>0.6</b>	0.41	nan	0.37
unclear	<b>0.14</b>	-0.01	nan	0.07
heterosexual	<b>0.45</b>	0.11	nan	0.39
not-referring	0.52	0.35	nan	<b>0.54</b>
other	0.29	0.26	nan	<b>0.54</b>
homosexual	0.79	<b>0.89</b>	nan	0.67
bisexual	0.46	<b>0.64</b>	nan	0.33

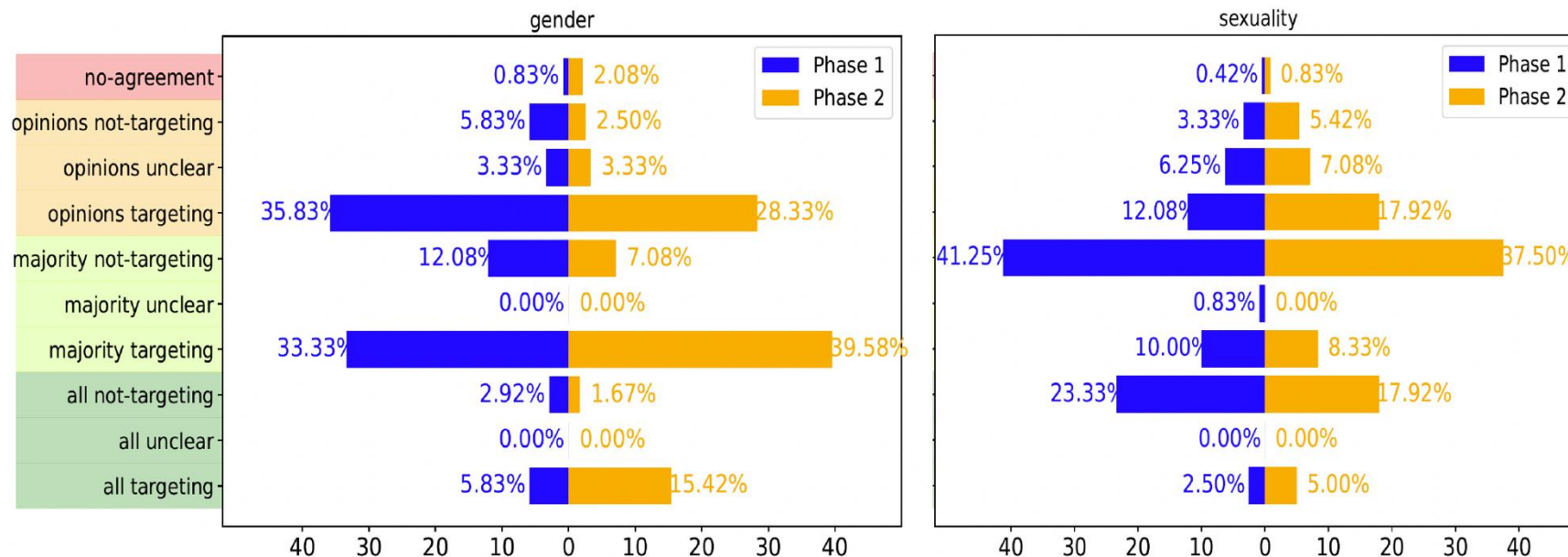
Gender Labels: non-cisgender annotators (G) as the group targeted

Sexuality Labels: non-heterosexual annotators (S) as the group targeted

When semantics increased inter-annotator agreement, it was because annotations from non-target groups **aligned more closely** with those from the target group.

# Results

## 3. How does the identification of hate target groups change after introducing semantics?



*Each text may be targeting, not targeting, or unclear as decided by all, the majority, or at least two annotators (opinions).*

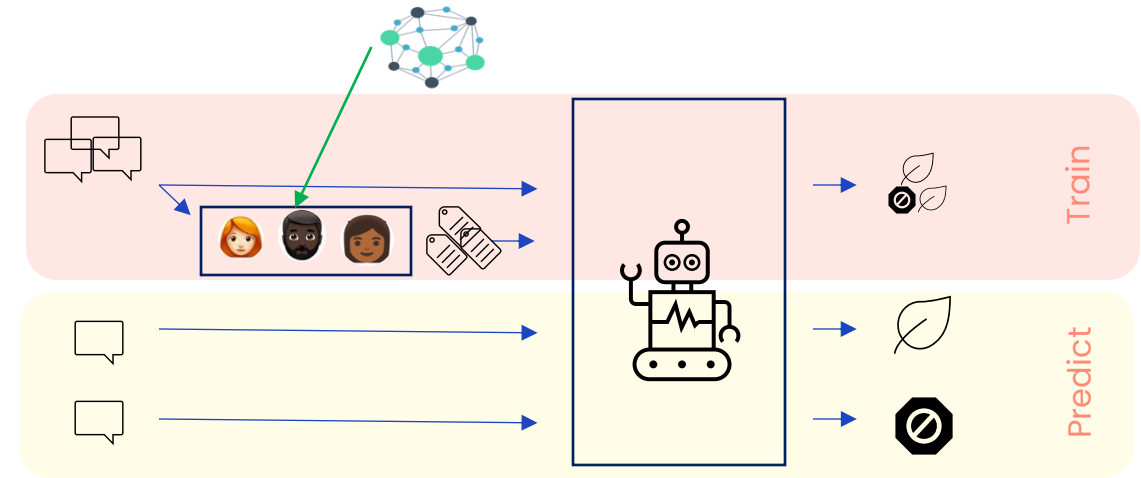
Semantics helps to **resolve hard-to-annotate text** (i.e., posts with a lower agreement or initially not deemed as pointing to hate-speech target groups).

# Conclusion

Background semantics is key for enhancing

- Inter-annotator agreement
- Alignment with annotators from the targeted group
- Comprehension and knowledge acquisition in all groups of annotators

## Incorporating Semantics in Hate Speech Annotation

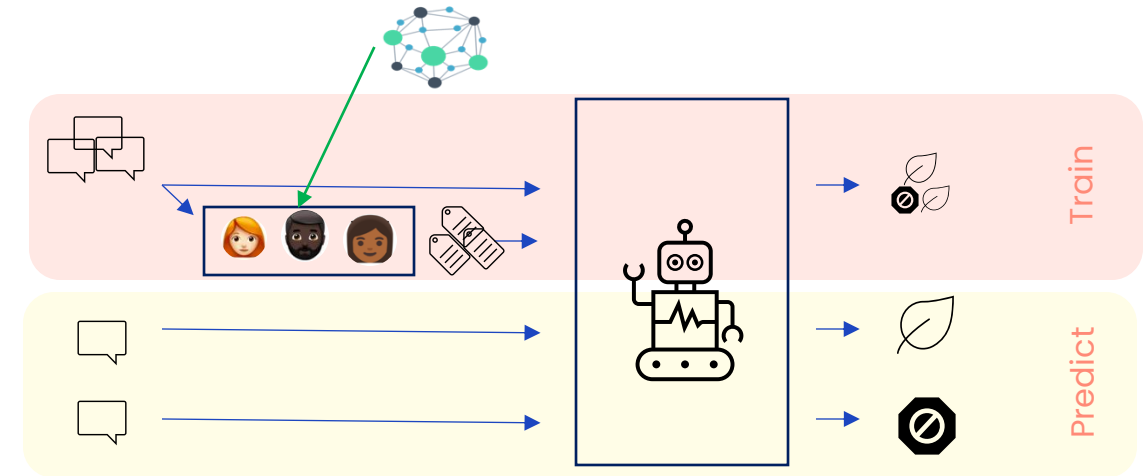


# Conclusion

Background semantics is key for enhancing

- Inter-annotator agreement
- Alignment with annotators from the targeted group
- Comprehension and knowledge acquisition in all groups of annotators

## Incorporating Semantics in Hate Speech Annotation



## Limitations

- Semantics being less effective when hate speech is implicit
- Testing impact on ML performance

# Semantics helps annotators to process hateful terminology often used to target individuals or groups in social media

*Enhancing Hate Speech Annotations with Background Semantics*

*Reyero Lobo, P., Daga, E., Alani, H., & Fernandez, M.*

## Contact

[\*paula.reyero-lobo@open.ac.uk\*](mailto:paula.reyero-lobo@open.ac.uk)



Funded by  
the European Union



Data Repository



Code Repository

