

Knowledge-Grounded Target Group Language Recognition in Hate Speech

Paula Reyero Lobo

Knowledge Media Institute

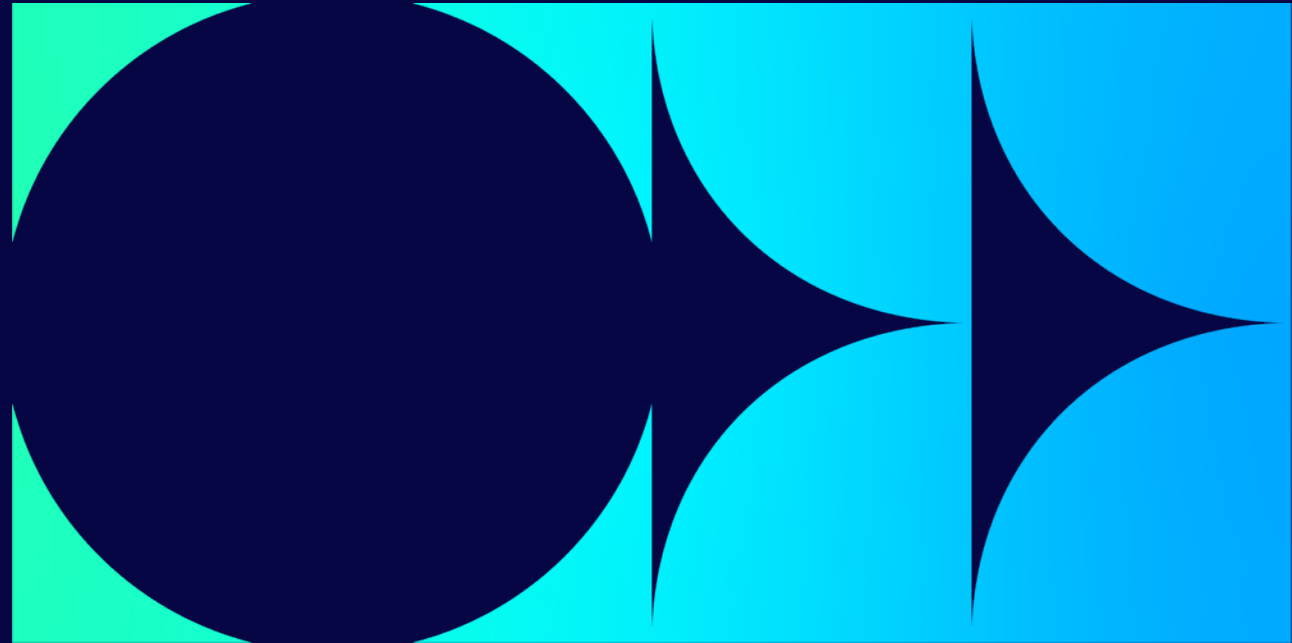
20-22 September 2023

Advised by:

Enrico Daga, Harith Alani, Miriam Fernandez



Funded by
the European Union



Introduction

Does this post mention or is about...

gender

sexuality

race...?



It's been 3 months and these words have only proven to be true. Suck it Butch who can't Deadlift.



Is this that bimbo twat with fish lips and hideous full sleeve tats?



Introduction

Does this post mention or is about...

gender
sexuality
race...?

Hate speech contains **nuanced, specialized language** to attack frequently target groups.

It's been 3 months and these words have only proven to be true. Suck it Butch who can't Deadlift.

Is this that bimbo twat with fish lips and hideous full sleeve tats?

Introduction

Does this post mention or is about...

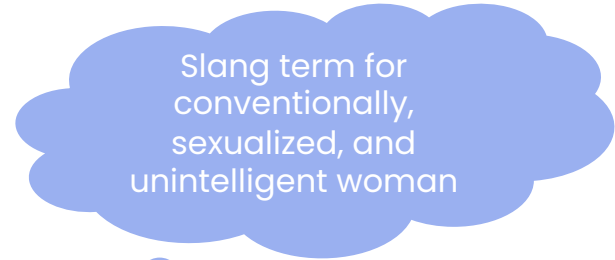
gender
sexuality
race...?



It's been 3 months and these words have only proven to be true. Suck it **Butch** who can't Deadlift.



Is this that **bimbo** twat with **fish lips** and hideous full sleeve tats?



Problem statement

Does this post mention or is about **gender or sexuality**?

Post 1

It's been 3 months and these words have only proven to be true. Suck it **Butch** who can't Deadlift.

Label



Sexuality



None

Post 2

Is this that **bimbo** twat with **fish lips** and hideous full sleeve tats?

Label



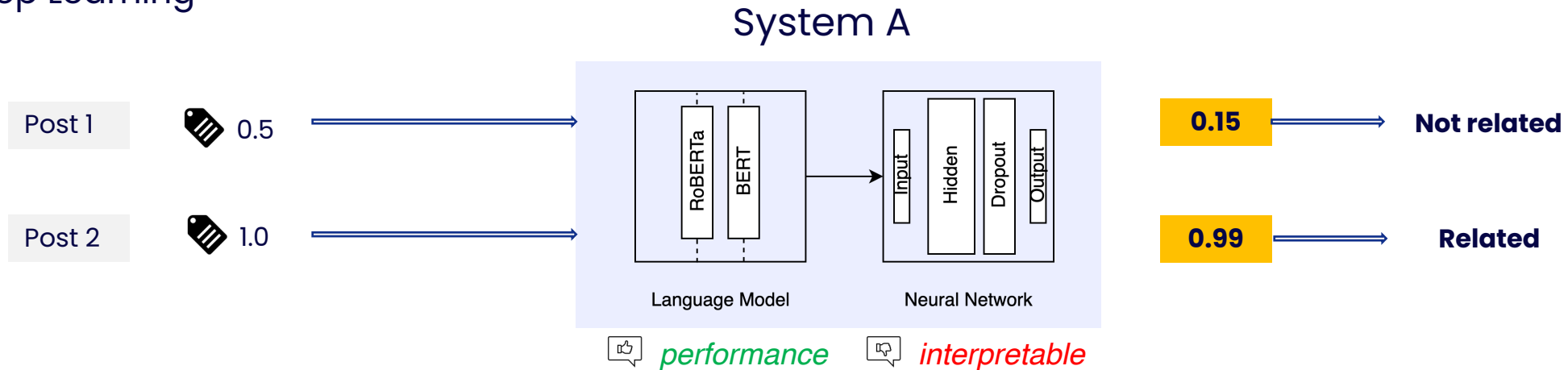
Gender

The label bias in training datasets calls for the need of **interpretability** when building automated detection systems.

Related Work

Does this post mention or is about **gender or sexuality**?

▪ Deep Learning



Post 1

It's been 3 months and these words have only proven to be true. Suck it Butch who can't Deadlift.

Label



Post 2

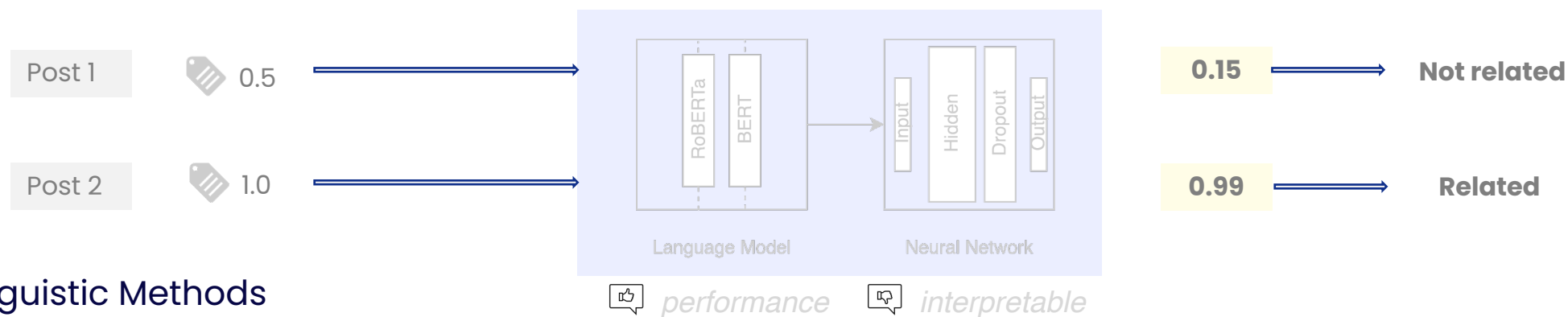
Is this that bimbo twat with fish lips and hideous full sleeve tattts?



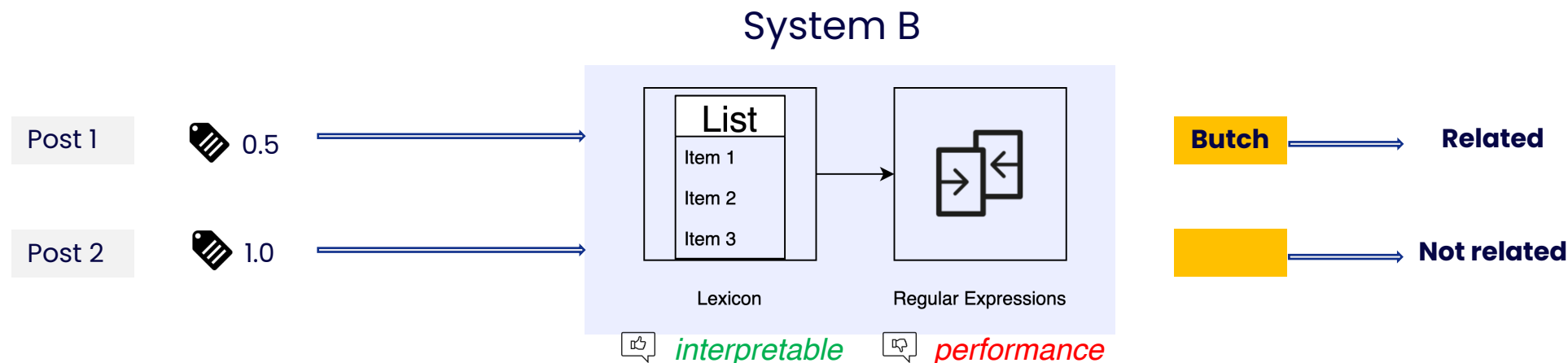
Related Work

Does this post mention or is about **gender or sexuality**?

Deep Learning



Linguistic Methods



Post 1

It's been 3 months and these words have only proven to be true. Suck it Butch who can't Deadlift.

Label



Post 2

Is this that bimbo twat with fish lips and hideous full sleeve tattts?

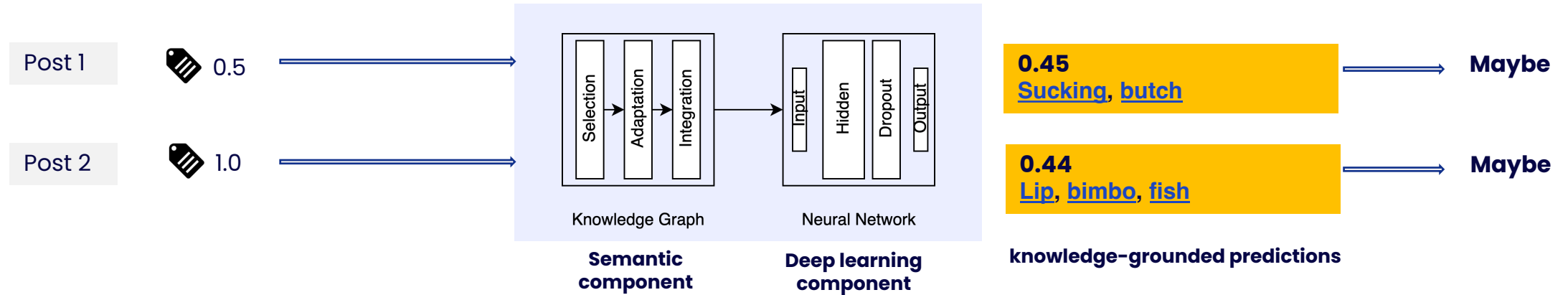


Our approach

Does this post mention or is about **gender or sexuality**?

▪ Hybrid Learning

Our System



Our approach

Does this post mention or is about **gender or sexuality**?

▪ Hybrid Learning

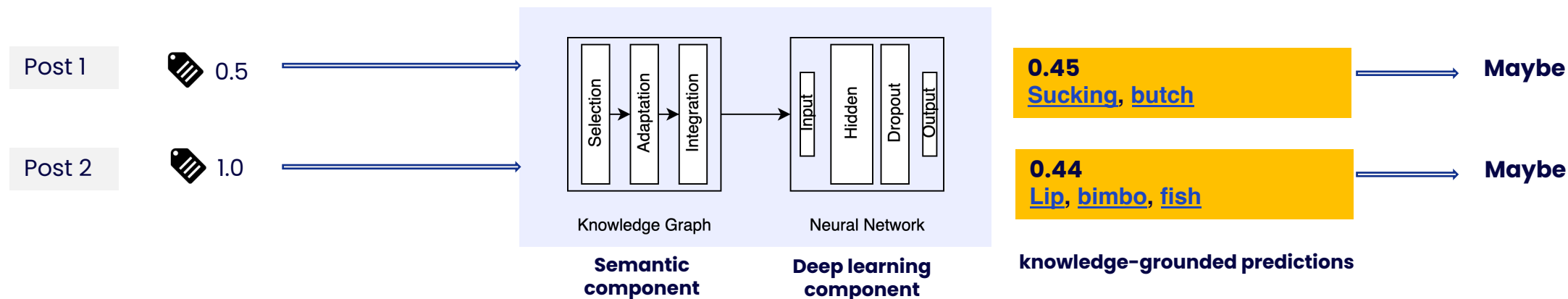
Post 1

It's been 3 months and these words have only proven to be true. Suck it Butch who can't Deadlift.

Post 2

Is this that bimbo twat with fish lips and hideous full sleeve tats?

Our System



Instance: **butch**

Term IRI: http://purl.obolibrary.org/obo/GSSO_000338

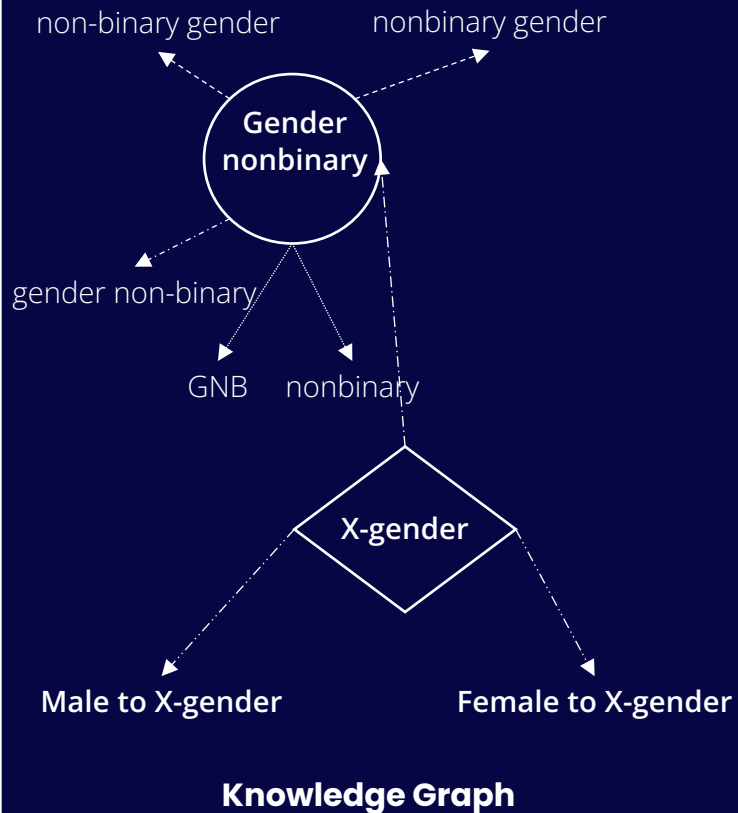
Definition: Traditionally, a lesbian who appears "masculine" or acts in a "masculine" manner. However, the term is often used to describe a wide variety of gender expressions and identities adjacent to masculinity or "butchility".

Annotations

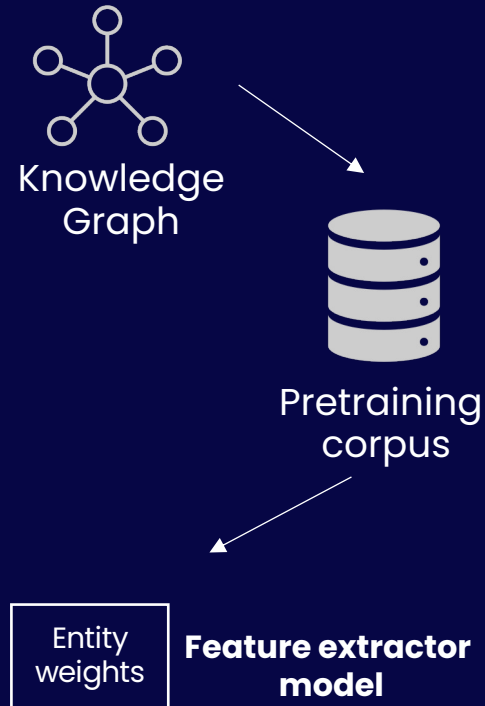
- **comment:** In the first half of the 20th century, "butch" was used as a gay male slang term referring to a masculine gay man.

Hybrid Approach

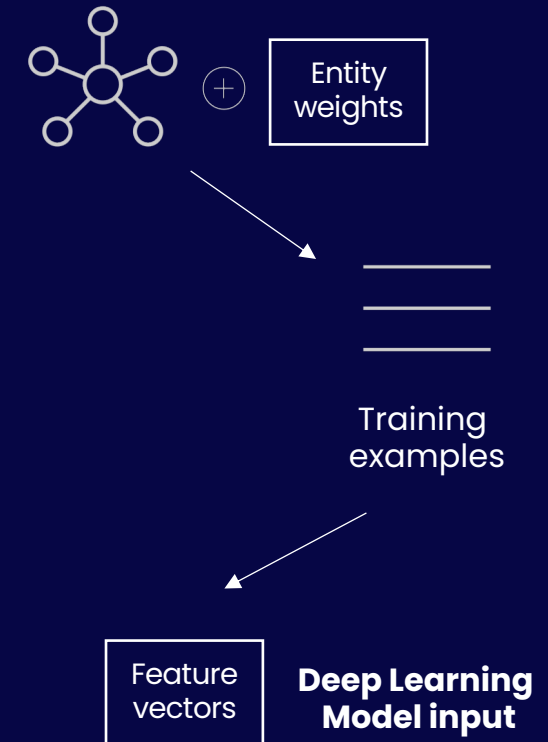
1. Selection



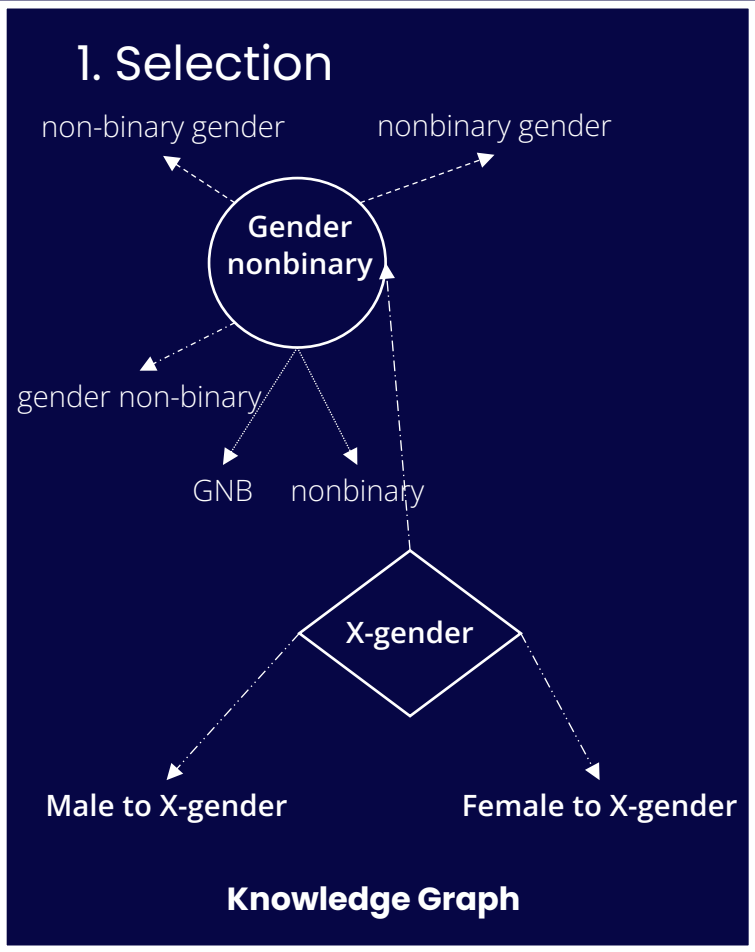
2. Adaptation



3. Integration



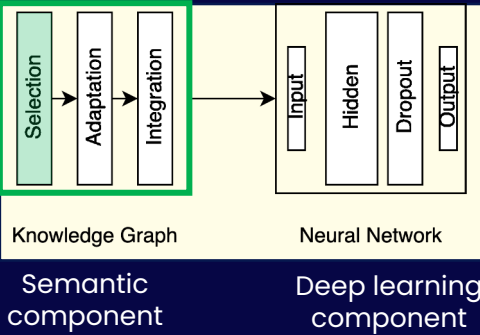
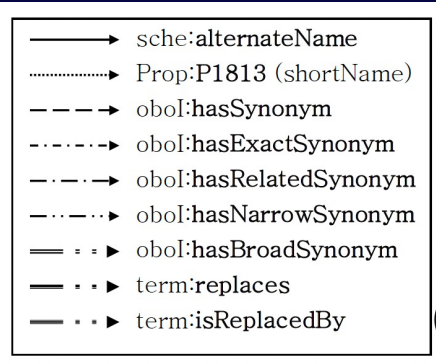
Hybrid Approach



➤ Knowledge graph



Gender, Sex, and Sexual Orientation
Ontology



Hybrid Approach

➤ Pretraining data



Jigsaw Toxicity Sample

2. Adaptation



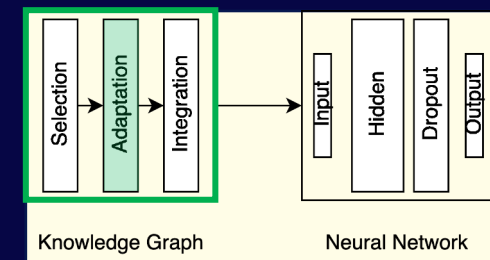
Knowledge Graph



Pretraining corpus

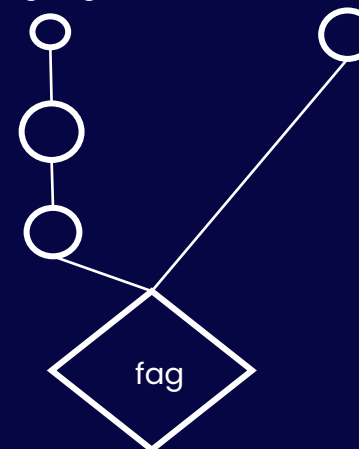
Entity weights

Feature extractor model

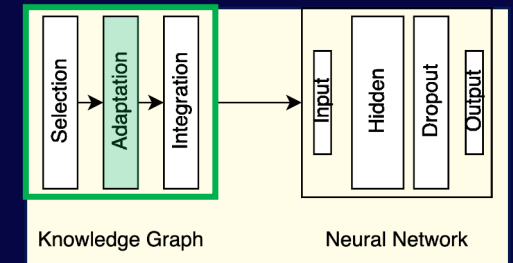
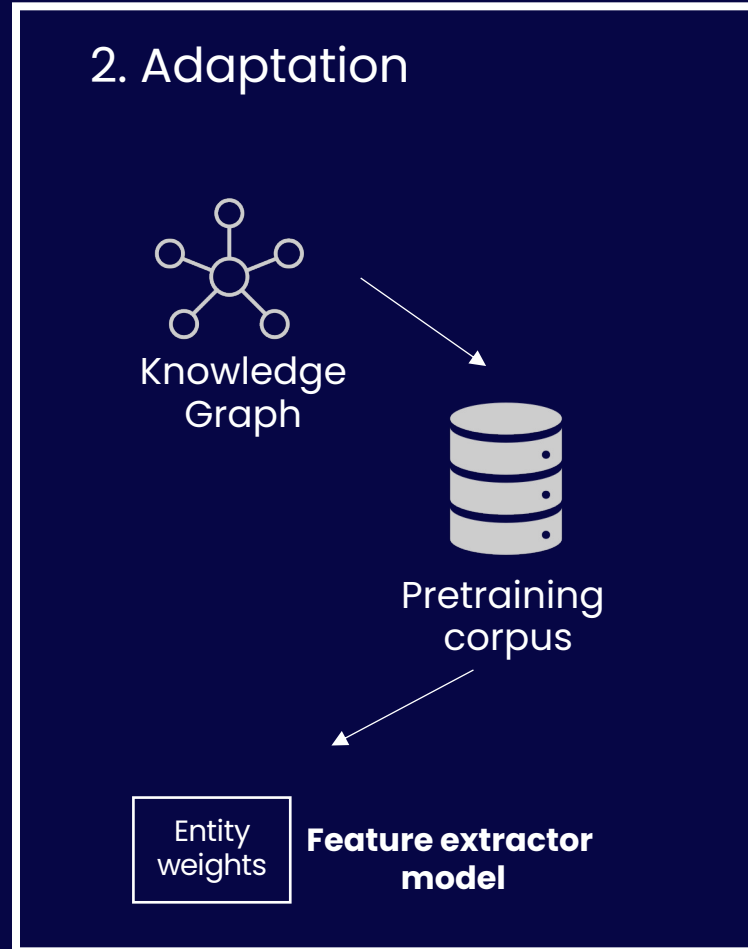


Figurative language

Homophobic slur



Hybrid Approach



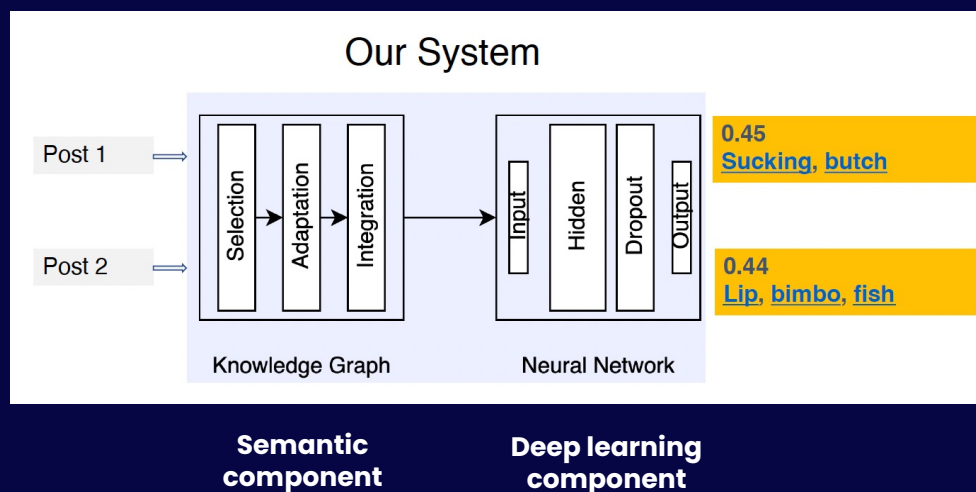
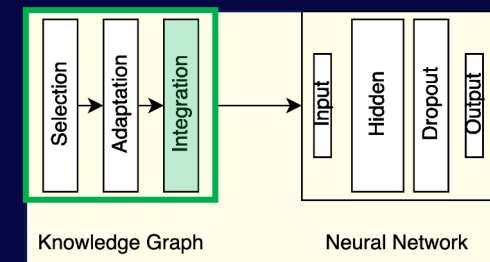
➤ Entity weighting*

Frequency-based
DocF

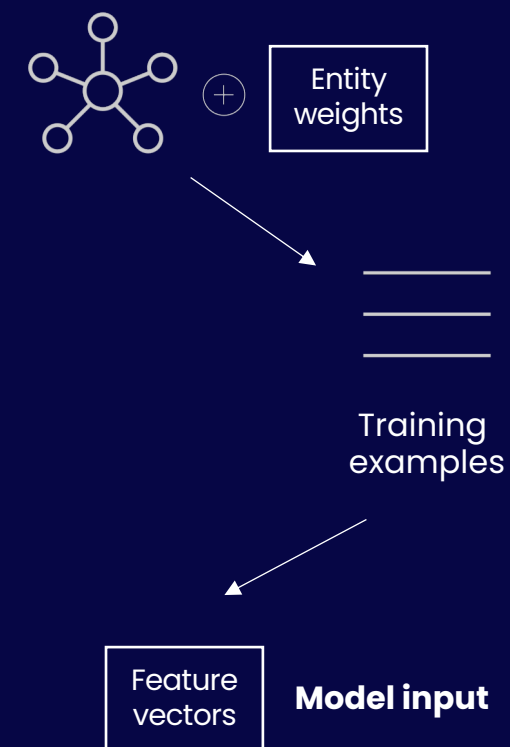
Machine Learning coefficients
Logistic Regression
Multinomial Naïve Bayes

* Hierarchical entity expansion: an entity (e_i) appears in a text if e_i , any of its subclasses, or any of its instances appears in the text

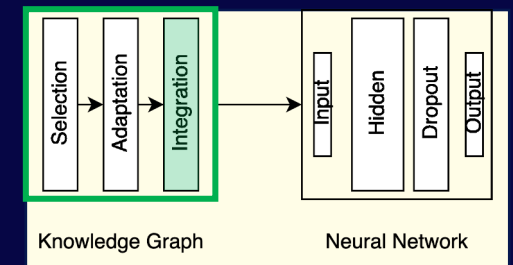
Hybrid Approach



3. Integration



Hybrid Approach



➤ Training data

Measuring Hate Speech

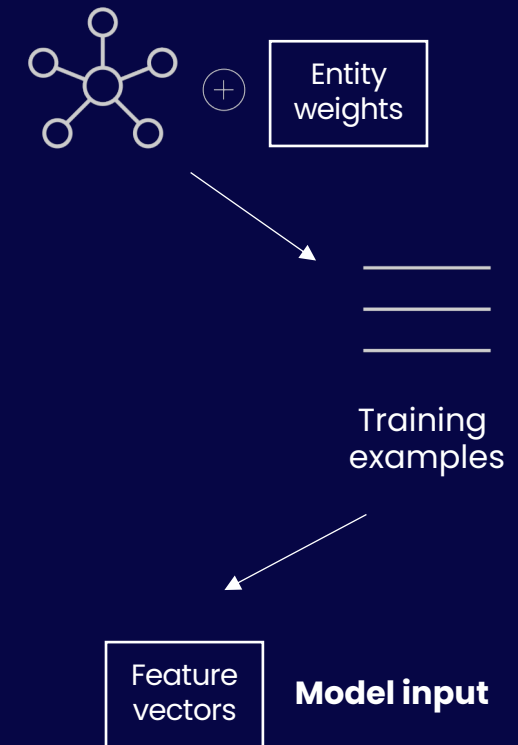
➤ Evaluation data

Gab Hate Corpus
HateXplain
Xtremespeech

➤ Baselines

RoBERTa base (System A)
Toxic Debias (System B)

3. Integration



Sachdeva, Patrik, et al. "The Measuring Hate Speech Corpus: Leveraging Rasch Measurement Theory for Data Perspectivism." In *Proceedings of the 1st Workshop on Perspectivist Approaches to NLP*. 2022.

Sachdeva, Patrik, et al. "Targeted Identity Group Prediction in Hate Speech Corpora.." In: *Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)*. 2022.

Zhou, Xuhui, et al. "Challenges in Automated Debiasing for Toxic Language Detection." *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*. 2021.

Results

Quantitative analysis

	Model	Gab Hate Corpus				XtremeSpeech _{English}				HateXplain			
		Accuracy	<i>F1</i>	ROC AUC	PR AUC	Accuracy	<i>F1</i>	ROC AUC	PR AUC	Accuracy	<i>F1</i>	ROC AUC	PR AUC
Lexicon ¹	Toxic Debias	91.81	58.82	74.82	40.20	94.01	52.41	72.96	31.12	84.43	67.36	79.01	52.66
Hybrid	HybridDocF _h	91.30	51.15	84.52	54.55	93.97	53.91	87.02	47.05	79.45	43.55	78.37	55.96
	HybridLR _h	90.64	62.42	89.30	64.38	90.79	49.27	88.36	50.79	83.48	67.72	88.15	68.35
	HybridMultiNB _h	89.36	61.11	90.13	68.24	90.38	47.74	87.26	51.80	85.63	73.57	91.38	78.37
PLM ²	RoBERTa _{base}	88.85	61.55	93.06	70.32	92.99	57.67	93.67	57.38	89.91	80.22	95.60	86.46

Our hybrid approach outperforms existing interpretable approaches and is, to some extent, competitive with the use of pre-trained language models (PLMs).

Results

Qualitative analysis

We identify gender and sexuality related language in samples of human and model (dis)agreement.

Does this post mention or is about **gender or sexuality**?

I don't think I am being gullible ,[...].Not to mention **homosexuals** and the mentally ill .If it is an international racist insult it falls flat .



Not related



0.83

gay, homosexuality, insult

Our hybrid approach is interpretable in that it identifies features that contribute to predict the class.

agree



woman, man, LGBTQ, transgender person, Black man, gay, homo, heterosexual, homosexuality, feminist, lesbian, queer sexual orientation, same-gender marriage, pederasty, sexuality, homophobia, queer identity, *sex, feminism, marital partner, sex worker, fag, faggot, masculism, pussy, thot, rape, menstruation, bitch*

disagree



woman, man, LGBT, woman of color, **gay**, heterosexual, homosexuality, feminist, lesbian, queer sexual orientation, transgender, same-gender marriage, womanism, sexuality, homophobia, queer identity, abusive person, sex, *partner, faggot, semen, pussy, bitch*



feminist, homophobia, sex, sex work client, *marital partner, parent, sex worker, fag, faggot, rapist, female gender role, pussy, abortion, vagina, thot, penis, rape, domestic violence, she, bimbo, sexual abstinence, cunt, bitch, whore, slut, fuck, Mrs., rainbow flag*

In our case, the language related to gender and sexuality.

What's wrong with that? It's History you know. Mary Beard told me so and she has a **vagina** from Oxford.



Related to gender



0.49

beard, vagina, she

Results

Qualitative analysis

Our error analysis finds that 97% of the false positives were **annotation errors**.

The entities in red were mentioned in texts that annotators deemed as “not related”.



woman, man, LGBTQ, transgender person, Black man, gay, homo, heterosexual, homosexuality, feminist, lesbian, queer sexual orientation, same-gender marriage, pederasty, sexuality, homophobia, queer identity, *sex, feminism, marital partner, sex worker, fag, faggot, masculism, pussy, thot, rape, menstruation, bitch*



woman, man, LGBT, woman of color, **gay**, heterosexual, homosexuality, feminist, lesbian, queer sexual orientation, transgender, same-gender marriage, womanism, sexuality, homophobia, queer identity, abusive person, sex, *partner, faggot, semen, pussy, bitch*



feminist, homophobia, sex, sex work client, *marital partner, parent, sex worker, fag, faggot, rapist, female gender role, pussy, abortion, vagina, thot, penis, rape, domestic violence, she, bimbo, sexual abstinence, cunt, bitch, whore, slut, fuck, Mrs., rainbow flag*

Does this post mention or is about **gender or sexuality**?

I don't think I am being gullible, [...].Not to mention **homosexuals** and the mentally ill .If it is an international racist insult it falls flat .



Not related



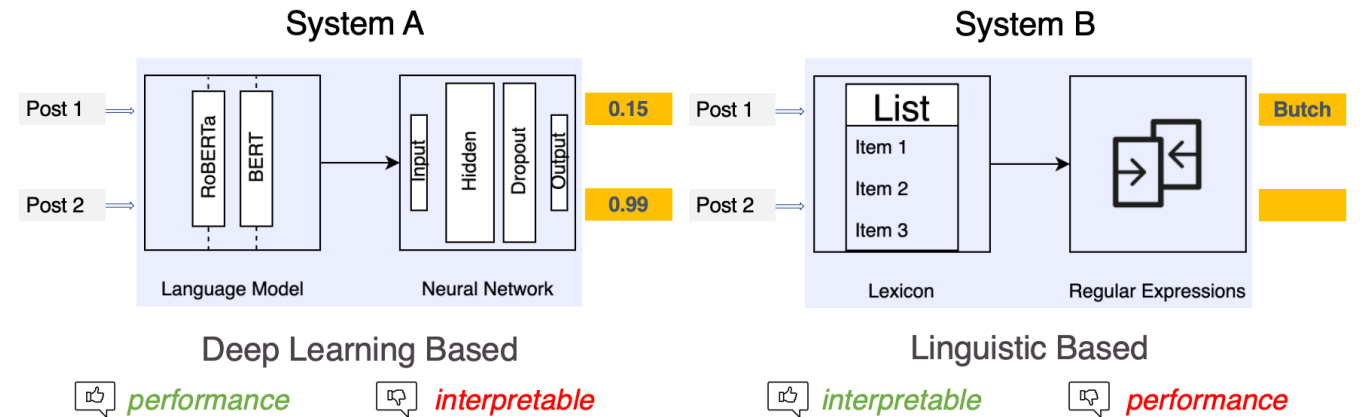
0.83

gay, homosexuality, insult

Conclusion

Hybrid learning is key for enhancing

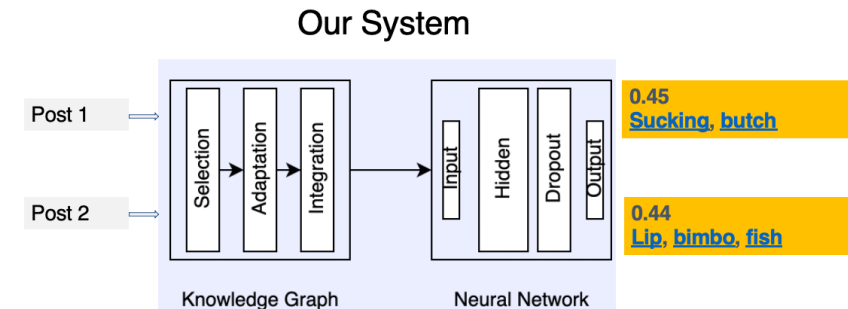
- Interpretability
- Robustness
- Handling annotation errors



Limitations

- KG are human-intensive
- Moving from use case to a general case scenario

HYBRID LEARNING



Open questions

To what extent a hybrid architecture can...

- *be an extrinsic interpretation of deep learning/opaque models?*
- *prevent or correct data annotation errors?*
- *audit bias in hate speech detection systems?*

Combining a KG with deep learning provides interpretability in the detection of target groups in hate speech

Knowledge-Grounded Target Group Language Recognition in Hate Speech

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

Contact

paula.reyero-lobo@open.ac.uk



www.linkedin.com/in/paula-reyero-lobo-116449170



<https://orcid.org/0000-0001-5238-4550>



Funded by
the European Union

To what extent a hybrid architecture can...

- *interpret deep learning/opaque models?*
- *prevent or correct data annotation errors?*
- *audit bias in hate speech detection systems?*

