# FEATURES OR SPURIOUS ARTIFACTS?
## Data-centric baselines for fair and robust hate speech detection

*Alan Ramponi, Sara Tonelli – Fondazione Bruno Kessler, Italy*

## LEXICAL ARTIFACTS

⚠️ *Hate speech detection models overly rely on lexical artifacts, affecting **fairness** & **robustness***

## LEXICAL ARTIFACTS ACROSS PLATFORMS

Hate speech corpora from:

lang: **English**
mod: **written**

Reddit  Gab  Twitter  Stormfront

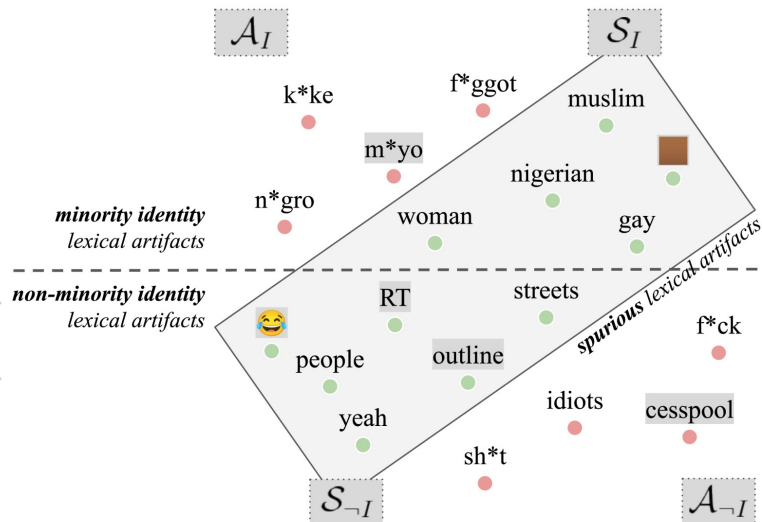**Cross-distribution PMI** & artifacts annotation

## ARTIFACTS STATEMENT

**Documentation** *of lexical artifacts to contribute to diagnosis & mitigation of pre-existing biases:*

**Template:**
I.   TOP LEXICAL ARTIFACTS
II.  CLASS DEFINITIONS
III. METHODS & RESOURCES

## ARE LEXICAL ARTIFACTS *ALL THE SAME*?



$\mathcal{A}_I$   $\mathcal{S}_I$
k*ke   f*ggot   muslim
m*yo   nigerian
*minority identity lexical artifacts*   n*gro   woman   gay
*non-minority identity lexical artifacts*   RT   streets   *spurious lexical artifacts*
people   outline   f*ck
yeah   idiots   cesspool
sh*t
$\mathcal{S}_{\neg I}$   $\mathcal{A}_{\neg I}$

## DATA-CENTRIC LEXICAL DEBIASING

Use <u>identity</u> ($\mathcal{S}_I$) and <u>non-identity</u> ($\mathcal{S}_{\neg I}$) **spurious artifacts** for debiasing models

- **Fairness:** MASK $\mathcal{S}_I$ strong baseline
- **Robustness:** Hard to achieve with artifacts or generic solutions only

*Fairness & robustness **should be studied together** in future research*

## RESOURCES

🌐 github.com/dhfbk/hate-speech-artifacts