

Equally Safe Online?

46%

of women have received abuse

50%

of Black and minority ethnic women

(Glitch & EVAW, 2020)

23%

of women comfortable expressing political views online

(Stephens et al., 2024)

Automatic Identification and Classification of Misogynistic Language on Twitter

Maria Anzovino¹, Elisabetta Fersini^{1(✉)}, and Paolo Rosso²

Overview of EXIST 2023: sEXism Identification in Social NeTworks

Laura Plaza^{1,2(✉)}, Jorge Carrillo-de-Albornoz^{1,2}, Roser Morante¹,
Enrique Amigó¹, Julio Gonzalo¹, Damiano Spina², and Paolo Rosso³

SemEval-2019 Task 5: Multilingual Detection of Hate Speech Against Immigrants and Women in Twitter

Valerio Basile[◇] Cristina Bosco[◇] Elisabetta Fersini[♡]
Debora Nozza[♡] Viviana Patti[◇] Francisco Rangel^{♣♣}
Paolo Rosso[♣] Manuela Sanguinetti[◇]

Developing a Multilingual Annotated Corpus of Misogyny and Aggression

Shiladitya Bhattacharya¹, Siddharth Singh², Ritesh Kumar², Akanksha Bansal³,
Akash Bhagat², Yogesh Dawer², Purnima Lohini⁴, Atul K. Gite³

"Be nice to your wife! The restaurants are closed": Can Gender Stereotype Detection Improve Sexism Classification?

Patricia Chiril and Farah Benamara and Véronique Moriceau

SWSR: A Chinese dataset and lexicon for online sexism detection

Aiqi Jiang^{a,*}, Xiaohan Yang^b, Yang Liu^a, Arkaitz Zubiaga^a

“Computer scientists should stop thinking about online hate speech as something requiring **methods**, and start thinking about it as something that demands **solutions**. This change — treating hate speech less like a task and more like the **real-world problem** it is — would **orient CS research towards the concerns of other stakeholders**, and thus begin the **collaborative pursuit** toward a safe Internet.”

Parker & Ruths, 2020

Equally Safe Online

Gavin Abercrombie Heriot-Watt University



Engineering and
Physical Sciences
Research Council

GLITCH



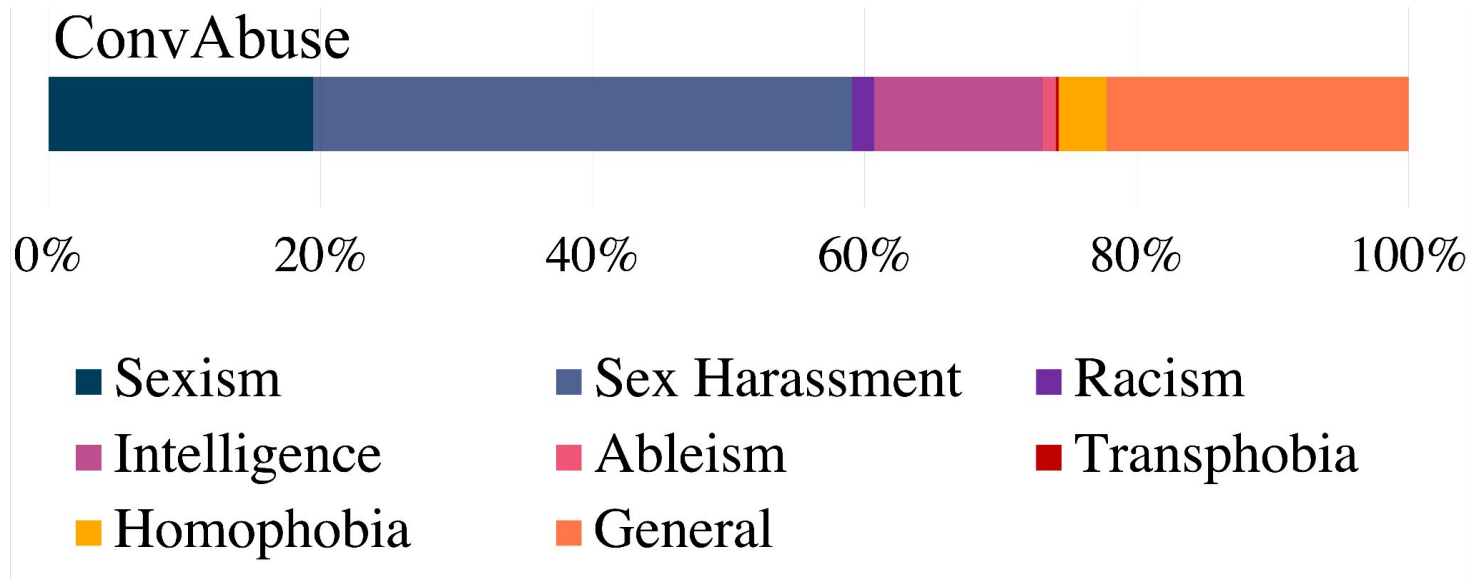
EmilyTest

Tackling Gender Based Violence in Education

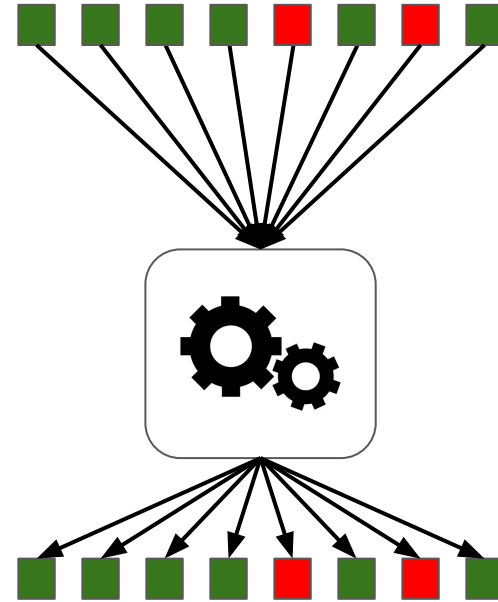
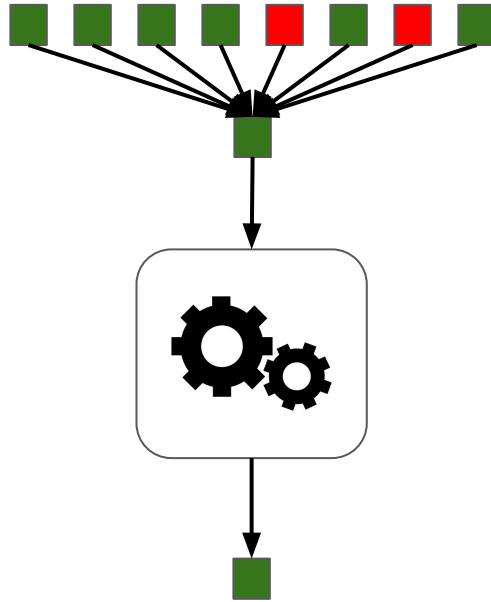
ConvAbuse: Data, Analysis, and Benchmarks for Nuanced Abuse Detection in Conversational AI

Amanda Cercas Curry¹ and Gavin Abercrombie¹ and Verena Rieser^{1,2}

EMNLP 2021



$$\alpha = 0.63$$



‘majority vote aggregation erases minoritized perspectives by choosing as the “right” label whatever is chosen by the majority of annotators, implicitly foregrounding dominant understandings and language ideologies.’

Blodgett, 2021

Say **Anything** with Boyfriend :)

Pron
Adv
Noun

Gimpel et al. 2011.

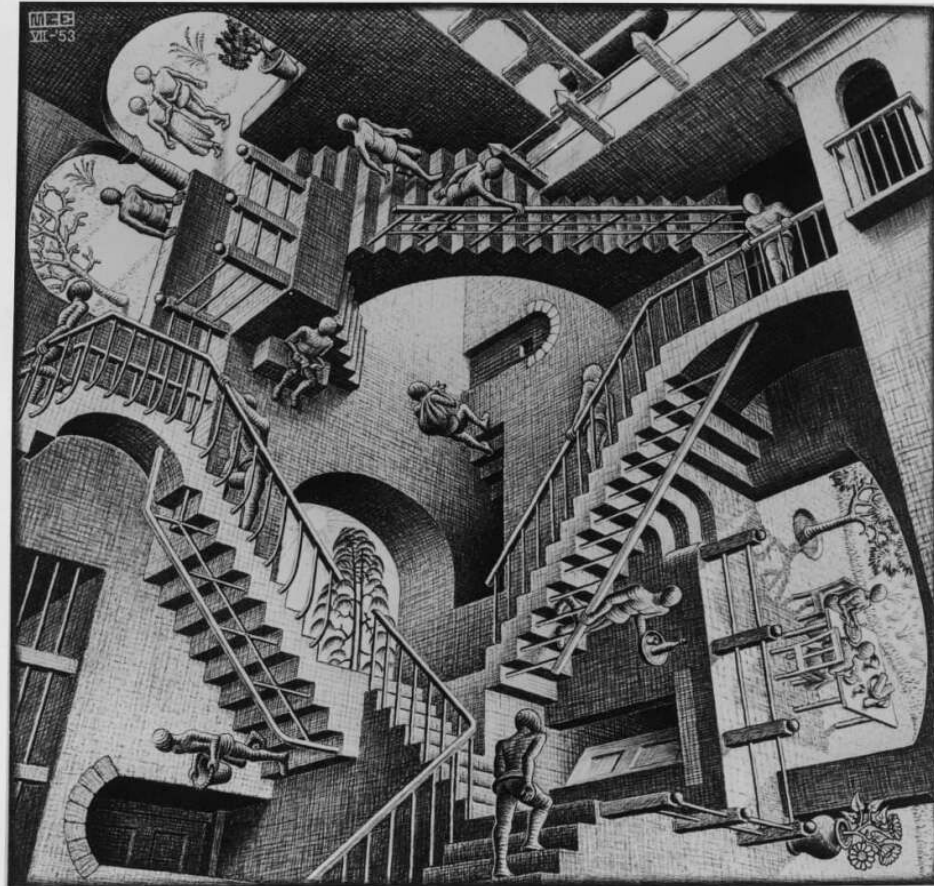
What is the background metal structure?



Ms COCO image id 393274, VQA 2.0 question id 393274004

- 1) trees
- 2) station
- 3) awning
- 4) platform
- 5) platform
- 6) platform
- 7) roof
- 8) shelter
- 9) train stop
- 10) awning

Goyal et al. 2017.



<https://nlperspectives.di.unito.it/>

LREC 2022 Marseille

ECAI 2023 Krakow

LREC-COLING 2024 Torino

SemEval-2023 Task 11: Learning With Disagreements (LeWiDi)

Elisa Leonardelli¹

Valerio Basile⁴

Verena Rieser²

Gavin Abercrombie²

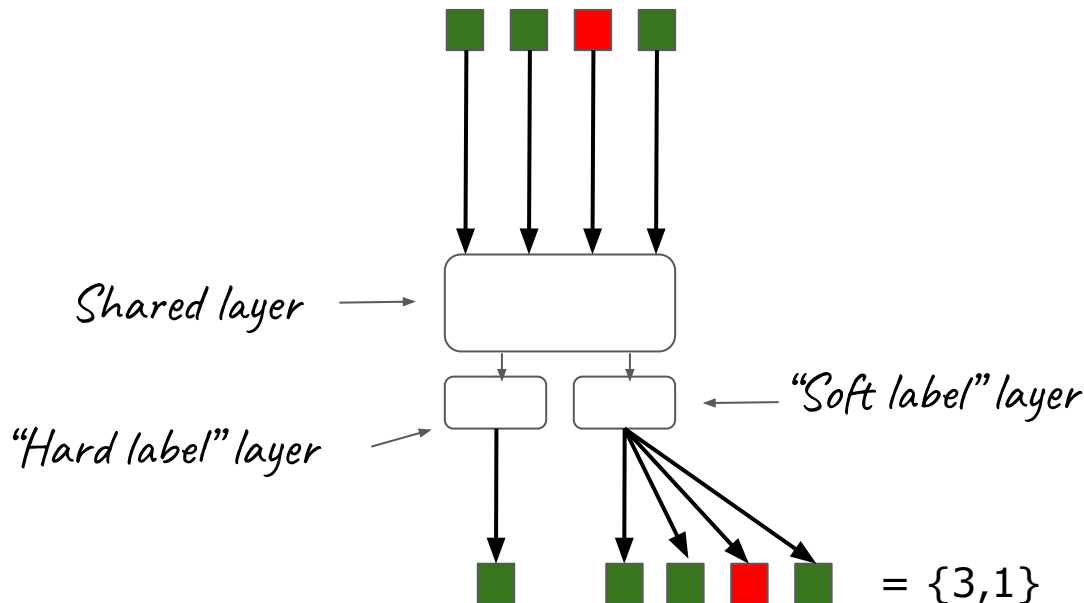
Tommaso Fornaciari⁵

Alexandra Uma

Dina Almanea³

Barbara Plank⁶

Massimo Poesio³



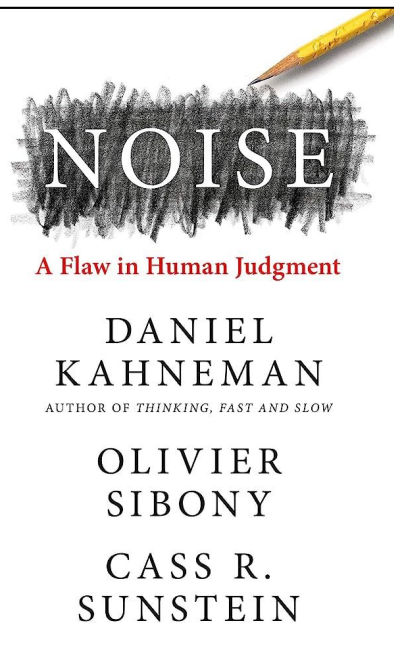
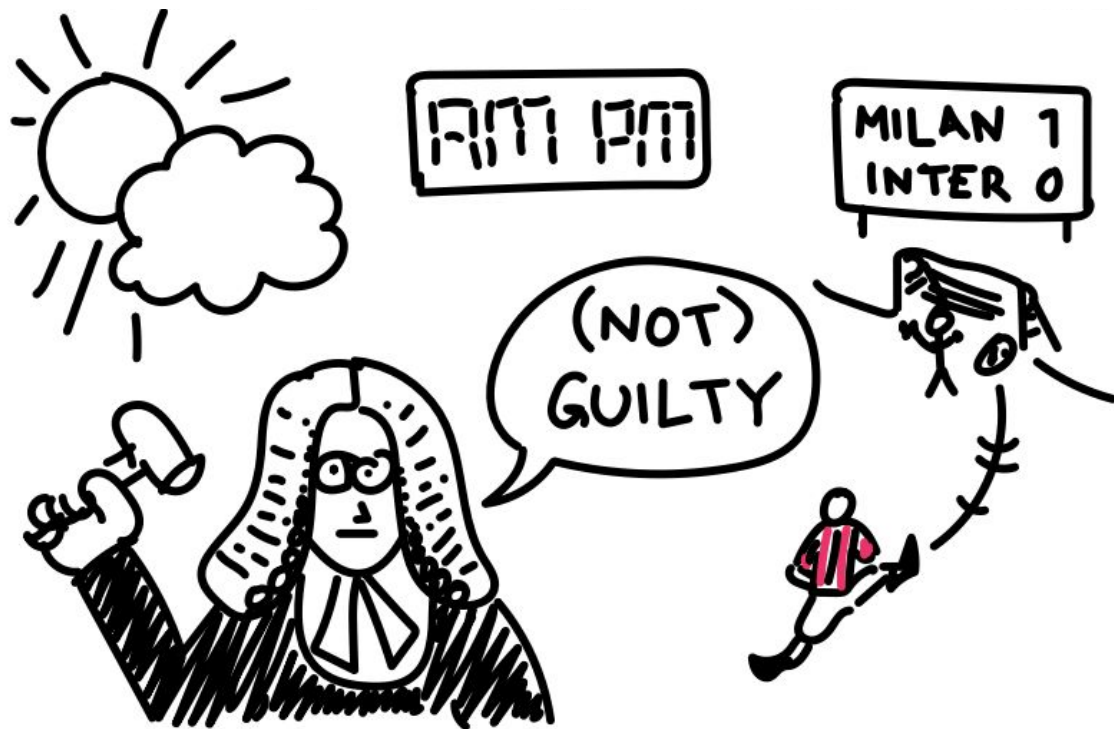
Temporal and Second Language Influence on Intra-Annotator Agreement and Stability in Hate Speech Labelling

Gavin Abercrombie
Heriot-Watt University

Dirk Hovy
Bocconi University

Vinodkumar Prabhakaran
Google Research

LAW 2023



Temporal and Second Language Influence on Intra-Annotator Agreement and Stability in Hate Speech Labelling

Gavin Abercrombie
Heriot-Watt University

Dirk Hovy
Bocconi University

Vinodkumar Prabhakaran
Google Research

74.5%
stability

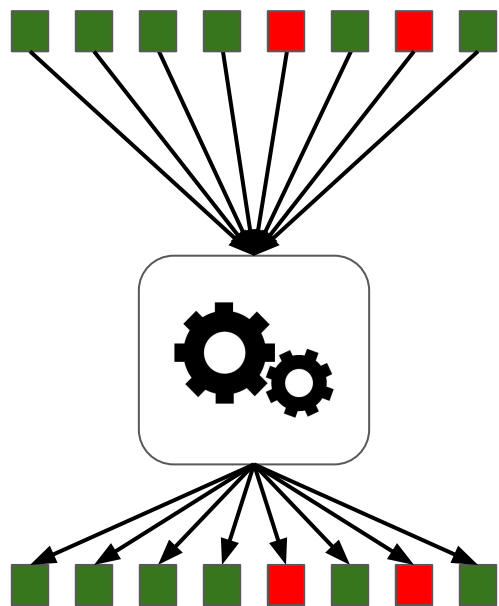
| | | Stability (temporal within annotator) | |
|--|-------------------|---|---|
| | | High <i>intra</i> | Low <i>intra</i> |
| Reliability (between annotators) | High <i>inter</i> | Straight-forward / Good quality | Systematic errors/ Value changes |
| | Low <i>inter</i> | Variable perspectives (high subjectivity) | Ambiguous or difficult / Poor quality |

Subjective *Isms*? On the Danger of Conflating Hate and Offence in Abusive Language Detection

Amanda Cercas Curry*

Gavin Abercrombie*

Zeerak Talat*



“When categories of abuse are described as subjective, we understand that there is no ground truth, and wider cultural norms do not impact what constitutes hate. Within the concept of *isms*, we argue that is the wrong approach and that these are culturally defined. That is, we argue that, for a stereotype or norm, there is a ground truth given by the cultural and temporal context a statement is made in.”

Task formulation
Annotator recruitment

Resources for Automated Identification of Online Gender-Based Violence: A Systematic Review

Gavin Abercrombie¹ and **Aiqi Jiang^{1,3}** and **Poppy Gerrard-Abbott^{4,5}**
and **Ioannis Konostas^{1,2}** and **Verena Rieser^{1*}**

R1. How is GBV characterised?

R2. Who is represented in annotation of the data?

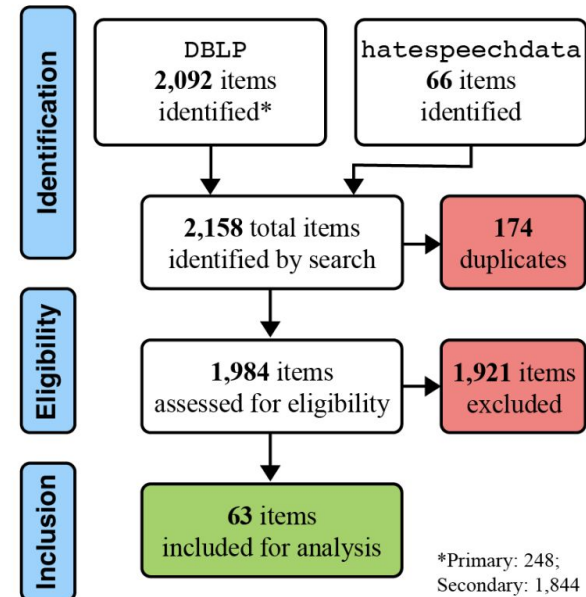
R3. From which platforms have the data been sourced?

R4. How has the data been sampled?

R5. Which languages are represented?

R6. During which time periods were the data created?

<https://github.com/HWU-NLP/GBV-Resources>



| | CONSULT | INCLUDE | COLLABORATE | OWN |
|------------------------------|---|--|---|---|
| <u>PARTICIPATION GOAL</u> | Why is participation needed? | | | |
| | To improve the user experience 80/80 | To better align AI with stakeholders' preferences and values 52/80 | To deliberate about system features 30/80 | To shape the system's scope and purpose 8/80 |
| <u>PARTICIPATION SCOPE</u> | What is on the table? | | | |
| | User interface of the system 80/80 | Underlying datasets (e.g., identification, curation, annotation) 8/80 | Overall design of system (e.g., task specification, model features) 8/80 | Whether and why the system should be built 4/80 |
| <u>FORM OF PARTICIPATION</u> | Who is involved? | | | |
| | Stakeholders recruited by the project team for discrete feedback 75/80 | Stakeholders recruited by the project team for domain expertise 47/80 | Stakeholders designated by the community collaborate in design 6/80 | Stakeholders designated by community play central role across project lifecycle 3/80 |
| | What form does stakeholder participation take? | | | |
| | Giving input on design ideas via questionnaires and interviews 68/80 | Group discussions with project team 49/80 | Ongoing collaborative prototyping and decision-making 18/80 | Reflexively deciding on the participatory approach 0/80 |

Delgado et al., 2023

Co-designing A Taxonomy of Online GBV

3 * workshops

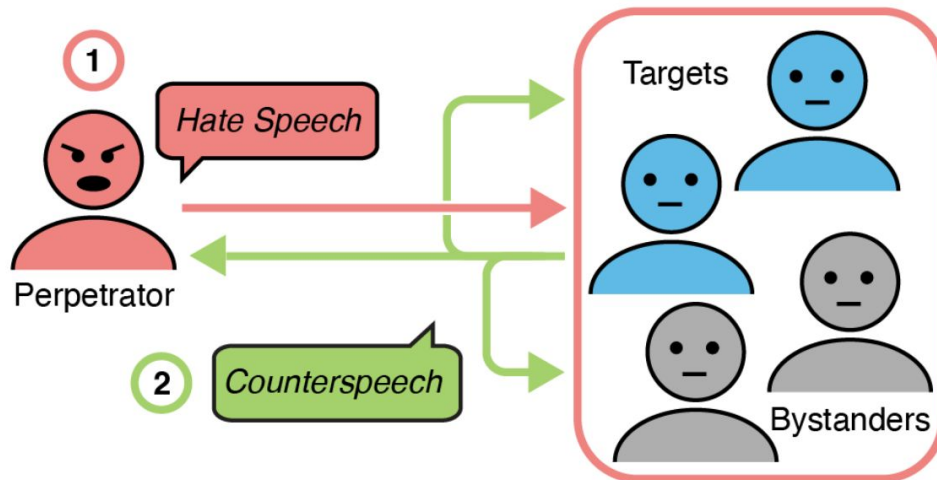
5 * focus groups

Severity cannot be judged from outside

Continuum of disempowerment to empowerment

Understanding Counterspeech for Online Harm Mitigation

Yi-Ling Chung¹ Gavin Abercrombie² Florence Enock¹
Jonathan Bright¹ Verena Rieser^{2*}



NLP for Counterspeech against Hate: A Survey and *How-To* Guide

Helena Bonaldi^{1,2} Yi-Ling Chung³ Gavin Abercrombie⁴ Marco Guerini¹

A Strategy Labelled Dataset of Counterspeech

Aashima Poudhar¹ and Ioannis Konstas^{1,2} and Gavin Abercrombie¹

Listen to stakeholders!

Gavin Abercrombie, Aiqi Jiang, Poppy Gerrard-abbott, Ioannis Konstas, and Verena Rieser. 2023. [Resources for Automated Identification of Online Gender-Based Violence: A Systematic Review](#). In *WOAH*. Association for Computational Linguistics.

Gavin Abercrombie, Dirk Hovy, and Vinodkumar Prabhakaran. 2023. [Temporal and Second Language Influence on Intra-Annotator Agreement and Stability in Hate Speech Labelling](#). In *Proceedings of LAW-XVII*. Association for Computational Linguistics.

Su Lin Blodgett. 2021. *Sociolinguistically Driven Approaches for Just Natural Language Processing*.

Helena Bonaldi, Yi-Ling Chung, Gavin Abercrombie, and Marco Guerini. 2024. [NLP for Counterspeech against Hate: A Survey and How-To Guide](#). In *Findings of NAACL*. ACL.

Amanda Cercas Curry, Gavin Abercrombie & Verena Rieser. 2021. Convabuse: Data, Analysis & Benchmarks for Nuanced Abuse Detection in Conversational AI. In *Proceedings of EMNLP*. ACL.

Amanda Cercas Curry, Gavin Abercrombie, and Zeerak Talat. 2024. [Subjective Isms? On the Danger of Conflating Hate and Offence in Abusive Language Detection](#). In *Proceedings of the 8th Workshop on Online Abuse and Harms (WOAH 2024)*, pages 275–282, Mexico City, Mexico. Association for Computational Linguistics.

Chung, Y. L., Abercrombie, G., Enock, F., Bright, J., & Rieser, V. (2023). Understanding counterspeech for online harm mitigation

Delgado, F., Yang, S., Madaio, M., & Yang, Q. 2023. The participatory turn in ai design: Theoretical foundations and the current state of practice. In *Proceedings of EEAMO*. ACM.

Gimpel et al. 2011. Part of Speech Tagging for Twitter: Annotation, features, and experiments.

Glitch UK and EVAW. 2020. The ripple effect: COVID-19 and the epidemic of online abuse.

Goyal et al. 2017. *Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering*.

Daniel Kahneman, Olivier Sibony & Cass R. Sunstein. 2021. *Noise. A Flaw in Human Judgement*.

Elisa Leonardelli, Gavin Abercrombie, Dina Almanea, Valerio Basile, Tommaso Fornaciari, Barbara Plank, Verena Rieser, Alexandra Uma, and Massimo Poesio. 2023. [SemEval-2023 Task 11: Learning with Disagreements \(LeWiDi\)](#). In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*. ACL.

Aashima Poudhar, Ioannis Konstas, and Gavin Abercrombie. 2024. [A Strategy Labelled Dataset of Counterspeech](#). In *Proceedings of WOAHS 2024*. ACL.

Stevens, F., Enock, F. E., Sippy, T., Bright, J., Cross, M., Johansson, P., ... & Margetts, H. Z. (2024). Women are less comfortable expressing opinions online than men and report heightened fears for safety: Surveying gender differences in experiences of online harms.