

7th Workshop on Online Abuse and Harms (WOAH 2023)

Toronto, Canada
13 July 2023

ISBN: 978-1-7138-8239-8

Printed from e-media with permission by:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571



Some format issues inherent in the e-media version may also appear in this print version.

Copyright© (2023) by the Association for Computational Linguistics
All rights reserved.

Printed with permission by Curran Associates, Inc. (2024)

For permission requests, please contact the Association for Computational Linguistics
at the address below.

Association for Computational Linguistics
209 N. Eighth Street
Stroudsburg, Pennsylvania 18360

Phone: 1-570-476-8006
Fax: 1-570-476-0860

acl@aclweb.org

Additional copies of this publication are available from:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571 USA
Phone: 845-758-0400
Fax: 845-758-2633
Email: curran@proceedings.com
Web: www.proceedings.com

ACL 2023

The 7th Workshop on Online Abuse and Harms (WOAH)

Proceedings of the Workshop

July 13, 2023

©2023 Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

Introduction

Digital technologies have brought many benefits for society, transforming how people connect, communicate and interact with each other. However, they have also enabled abusive and harmful content such as hate speech and harassment to reach large audiences, and for their negative effects to be amplified. The sheer amount of content shared online means that abuse and harm can only be tackled at scale with the help of computational tools. However, detecting and moderating online abuse and harms is a difficult task, with many technical, social, legal and ethical challenges. The Workshop on Online Harms and Abuse (WOAH) is the leading workshop dedicated to research addressing these challenges.

WOAH invites paper submissions from a wide range of fields, including natural language processing, machine learning, computational social sciences, law, politics, psychology, sociology and cultural studies. We explicitly encourage interdisciplinary submissions, technical as well as non-technical submissions, and submissions that focus on under-resourced languages. We also invite non-archival submissions for in progress work and reports from civil society to facilitate a meeting space between academic researchers and civil society.

This year marks the seventh edition of WOA, which will be co-located with ACL 2023 in Toronto, Canada. The special theme for this year’s edition is “subjectivity and disagreement in abusive language data”. Hate speech and other forms of abuse are highly subjective, in that there are diverse valid beliefs about what is or is not hateful or abusive. Different beliefs are informed by different social, cultural and legal norms. Through annotation, these beliefs are encoded in labelled datasets, which are then used to train and evaluate detection models. Therefore, subjectivity and disagreement are an essential aspect of research into online abuse and hate. By choosing this theme, we want to encourage submissions that examine, address or make use of this inherent subjectivity.

We received 55 submissions, of which 25 were accepted for presentation at the workshop. These papers will be presented at an in-person, where possible, poster session on the day of the workshop. Authors who are unable to attend in person will be able to give a virtual lightning talk describing their work. The workshop day will also include keynote talks from: Dirk Hovy, Milagros Miceli, Maarten Sap, Su Lin Blodgett, Vinodkumar Prabhakaran and Lauren Klein. Finally, we will close the day by inviting the keynote speakers to participate in a panel on the topic of subjectivity and disagreement.

We thank all our participants and reviewers for their work, and our sponsors for their support. We hope you enjoy this year’s WOA and the research published in these proceedings.

Paul, Yi-Ling, Debora, Aida, and Zeerak

Sponsors

WOAH is grateful for support from the following sponsors:

Diamond Tier



Gold Tier



Organizing Committee

Workshop Organiser

Yi-Ling Chung, The Alan Turing Institute
Aida Mostafazadeh Davani, Google Research
Debora Nozza, Bocconi University
Paul Röttger, University of Oxford
Zeeraq Talat, Independent Researcher

Program Committee

Emergency Reviewers

Gavin Abercrombie, Heriot Watt University
Greta Damo, Bocconi University
Lorenzo Lupo, Bocconi University

Program Committee

Syed Sarfaraz Akhtar, Apple Inc
Jisun An, Luddy School of Informatics, Computing, and Engineering, Indiana University Bloomington
Murali Raghu Babu Balusu, Georgia Institute of Technology
Francesco Barbieri, Snap Inc.
Valerio Basile, University of Turin
Thales Bertaglia, Maastricht University
Helena Bonaldi, Fondazione Bruno Kessler
Noah Broestl, University of Oxford, Google Research
Agostina Calabrese, The University of Edinburgh
Pedro Calais, UFMG, Brazil
Tommaso Caselli, Rijksuniversiteit Groningen
Amanda Cercas Curry, Bocconi University
Amit Das, Auburn University
Ona De Gibert, University of Helsinki
Daryna Dementieva, Technical University of Munich
Lucas Dixon, Google Research
Nemanja Djuric, Aurora Innovation
Hugo Jair Escalante, INAOE
Elisabetta Fersini, University of Milano-Bicocca
Bjørn Gambæk, Norwegian University of Science and Technology
Aitor García Pablos, Vicomtech
Sara Garza, FIME-UANL
Shlok Gilda, University of Florida
Lee Gillam, University of Surrey
Tonei Glavinic, Dangerous Speech Project
Darina Gold, Fraunhofer IIS
Marco Guerini, Fondazione Bruno Kessler
Udo Hahn, Friedrich-Schiller-Universität Jena
Christopher Homan, Rochester Institute of Technology
Muhammad Okky Ibrohim, University of Turin
Abhinav Jain, amazon.com
Mladen Karan, Queen Mary University
Mohammad Aflah Khan, IIT Delhi
Ian Kivlichan, Jigsaw, Google
Vasiliki Kougia, University of Vienna
Haewoon Kwak, Indiana University Bloomington
Sandra Kübler, Indiana University
Andrew Lee, University of Michigan
Els Lefever, LT3, Ghent University

Chuan-jie Lin, National Taiwan Ocean University
Nikola Ljubešić, Jožef Stefan Institute
Davide Locatelli, Technical University of Catalonia
Holly Lopez, Indiana University
Adrian Pastor Lopez Monroy, Mathematics Research Center CIMAT
Elizabeth Losh, William and Mary
Hongyin Luo, MIT
Sarah Masud, LCS2, IIITD
Puneet Mathur, University of Maryland College Park
Diana Maynard, University of Sheffield
Do June Min, University of Michigan
Manuel Montes, INAOE
Hamdy Mubarak, Qatar Computing Research Institute
Smruthi Mukund, Amazon
Preslav Nakov, Mohamed bin Zayed University of Artificial Intelligence
Isar Nejadgholi, National Research Council Canada
Brahmani Nutakki, Saarland University
Ali Omrani, University of Southern California
Matthias Orlikowski, Bielefeld University
Viviana Patti, University of Turin, Dipartimento di Informatica
Naiara Perez, Vicomtech
Matúš Pikuliak, Kempelen Institute of Intelligent Technologies
Flor Miriam Plaza-del-arco, Bocconi University
Michal Ptaszynski, Kitami Institute of Technology
Georg Rehm, DFKI
Bjorn Ross, University of Edinburgh
Molly Sauter, McGill University
Tyler Schnoebelen, Decoded AI
Alexandra Schofield, Harvey Mudd College
Indira Sen, GESIS
Caroline Sindere, Convocation Design + Research
Jeffrey Sorensen, Google Jigsaw
Gerasimos Spanakis, Maastricht University
Arjun Subramonian, University of California, Los Angeles
Afrin Sultana, University of Chittagong
Maite Taboada, Simon Fraser University
Radiathun Tasnia, University of Chittagong
James Thorne, KAIST AI
Zuoyu Tian, Indiana University
Dimitrios Tsarapatsanis, University of York
Avijit Vajpayee, Amazon
Francielle Vargas, University of São Paulo
Ruyuan Wan, University of Notre Dame
Jing Xu, Facebook AI
Qiangeng Yang, University of Florida
Seunghyun Yoon, Adobe Research
Aleš Završnik, Institute of criminology at the Faculty of Law Ljubljana
Torsten Zesch, Computational Linguistics, FernUniversität in Hagen
Yi Zheng, University of Edinburgh

Table of Contents

<i>Identity Construction in a Misogynist Incels Forum</i> Michael Yoder, Chloe Perry, David Brown, Kathleen Carley and Meredith Pruden	1
<i>DeTexD: A Benchmark Dataset for Delicate Text Detection</i> Artem Chernodub, Serhii Yavnyi, Oleksii Sliusarenko, Jade Razzaghi, Yichen Mo and Knar Hovakimyan	14
<i>Towards Safer Communities: Detecting Aggression and Offensive Language in Code-Mixed Tweets to Combat Cyberbullying</i> Nazia Nafis, Diptesh Kanojia, Naveen Saini and Rudra Murthy	29
<i>Towards Weakly-Supervised Hate Speech Classification Across Datasets</i> Yiping Jin, Leo Wanner, Vishakha Kadam and Alexander Shvets	42
<i>Respectful or Toxic? Using Zero-Shot Learning with Language Models to Detect Hate Speech</i> Flor Miriam Plaza-del-arco, Debora Nozza and Dirk Hovy	60
<i>Benchmarking Offensive and Abusive Language in Dutch Tweets</i> Tommaso Caselli and Hylke Van Der Veen	69
<i>Relationality and Offensive Speech: A Research Agenda</i> Razvan Amironesei and Mark Diaz	85
<i>Cross-Platform and Cross-Domain Abusive Language Detection with Supervised Contrastive Learning</i> Md Tawkat Islam Khondaker, Muhammad Abdul-mageed and Laks Lakshmanan, V.s.	96
<i>Aporophobia: An Overlooked Type of Toxic Language Targeting the Poor</i> Svetlana Kiritchenko, Georgina Curto Rex, Isar Nejadgholi and Kathleen C. Fraser	113
<i>Problematic Webpage Identification: A Trilogy of Hatespeech, Search Engines and GPT</i> Ojasvin Sood and Sandipan Dandapat	126
<i>Concept-Based Explanations to Test for False Causal Relationships Learned by Abusive Language Classifiers</i> Isar Nejadgholi, Svetlana Kiritchenko, Kathleen C. Fraser and Esmá Balkir	138
<i>Female Astronaut: Because sandwiches won't make themselves up there": Towards Multimodal misogyny detection in memes</i> Smriti Singh, Amritha Haridasan and Raymond Mooney	150
<i>Conversation Derailment Forecasting with Graph Convolutional Networks</i> Enas Altarawneh, Ameeta Agrawal, Michael Jenkin and Manos Papagelis	160
<i>Resources for Automated Identification of Online Gender-Based Violence: A Systematic Review</i> Gavin Abercrombie, Aiqi Jiang, Poppy Gerrard-abbott, Ioannis Konstas and Verena Rieser	170
<i>Evaluating the Effectiveness of Natural Language Inference for Hate Speech Detection in Languages with Limited Labeled Data</i> Janis Goldzycher, Moritz Preisig, Chantal Amrhein and Gerold Schneider	187
<i>HOMO-MEX: A Mexican Spanish Annotated Corpus for LGBT+phobia Detection on Twitter</i> Juan Vásquez, Scott Andersen, Gemma Bel-enguix, Helena Gómez-adorno and Sergio-luis Ojedatruera	202

<i>Factoring Hate Speech: A New Annotation Framework to Study Hate Speech in Social Media</i>	
Gal Ron, Effi Levi, Odelia Oshri and Shaul Shenhav	215
<i>Harmful Language Datasets: An Assessment of Robustness</i>	
Katerina Korre, John Pavlopoulos, Jeffrey Sorensen, Léo Laugier, Ion Androutsopoulos, Lucas Dixon and Alberto Barrón-cedeño	221
<i>Robust Hate Speech Detection in Social Media: A Cross-Dataset Empirical Evaluation</i>	
Dimosthenis Antypas and Jose Camacho-Collados	231