

Hate Speech Definition

Hate Speech

Abwerten von Personen oder Gruppen aufgrund von Identitätsmerkmalen wie ...

Breite Definition

- Sozialer Status/Bildung/ Einkommens-/ Berufsgruppe
- Politische Einstellung
- Körperlichen Merkmalen/ Aussehen

Enge Definition

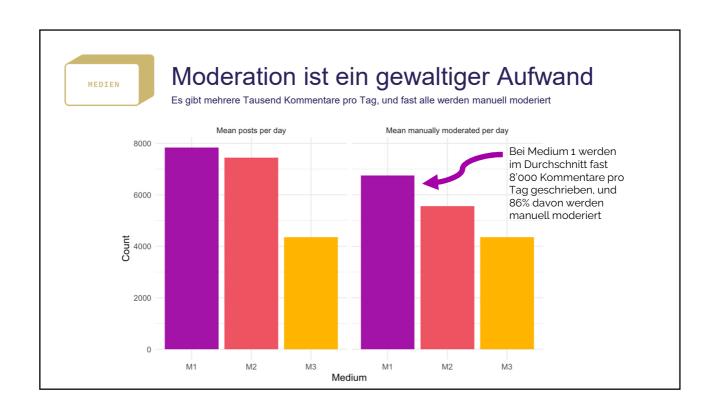
- Geschlecht
- Alter
- Sexualität
- Religion- Nationalität/ Hautfarbe/Herkunft
- Geistige/körperliche Beeinträchtigung

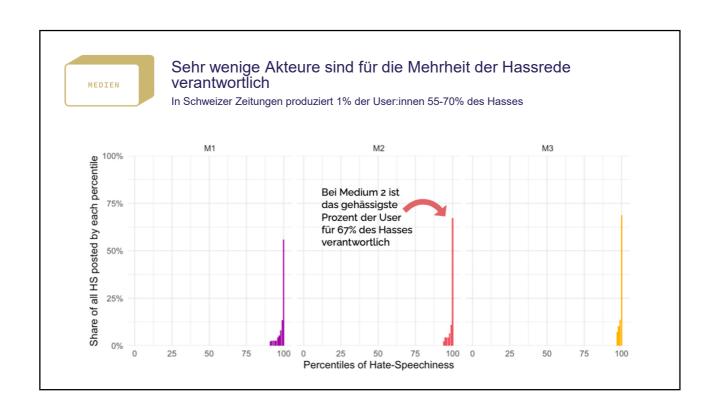
Toxische Sprache

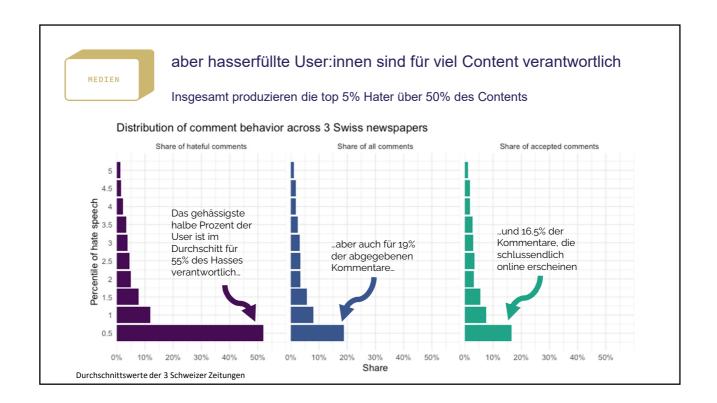
Verletzende, beleidigende oder vulgäre Sprache, die sich <u>nicht</u> auf Identitätsmerkmale bezieht

Und wir sehen...

- Unterrepräsentierte Gruppen ziehen sich weiter aus dem Diskurs zurück
- Nur eine kleine Minderheit beteiligt sich am Online-Diskurs oft die mit den polarisiertesten Meinungen
- Gesellschaftliche Spaltungen können somit verstärkt werden, mit Auswirkungen für die Demokratie
- Moderation bedeutet auch immer Zensur. Ein Gleichgewicht zu finden ist schwierig.









Counterspeech: ein Lösungsversuch

- Hate Speech mit Gegenrede statt Zensur bekämpfen
- Wie reagieren User:innen, wenn man auf Hate Speech antwortet?
- Wir haben verschiedene Strategien getestet



Empathy-based counterspeech can reduce racist hate speech in a social media field experiment

Dominik Hangartner^{a,b,1}, Gloria Gennaro^{a,b}o, Sary Alasiri^a, Nicholas Bahrich^a, Alexandra Bornhoft^a, Joseph Boucher^a, Buket Buse Demirci^a, Laurenz Derksen^{a,b}o, Aldo Hall^ao, Matthias Jochum^a, Maria Murias Munoz^a, Marc Richter^a, Franziska Vogel^a, Salomé Wittwer^a, Felix Wüthrich^a, Fabrizio Gilardi^c, and Karsten Donnay^co

