## Appendix C: Inclusion criteria and PRISMA diagram

*Query for dataset papers*

KEY ("social media" AND "dataset" AND ("NLP" OR "Natural Language Processing")) AND ("hate speech" OR "abus*" OR "offens*" OR "cyberbully*")) OR TITLE ("social media" AND "dataset" AND ("hate speech" OR "abus*" OR "offens*" OR "cyberbully*")) AND (LIMIT-TO (LANGUAGE, "English"))

*Inclusion criteria for dataset papers*

We applied the following inclusion criteria:
1) The paper presents a novel dataset for which annotation procedures are described.
2) The dataset is intended for training and testing algorithm(s) aimed at abuse detection.
3) The dataset is curated from one or more widely used social media platforms.
4) The dataset is in English.[1]
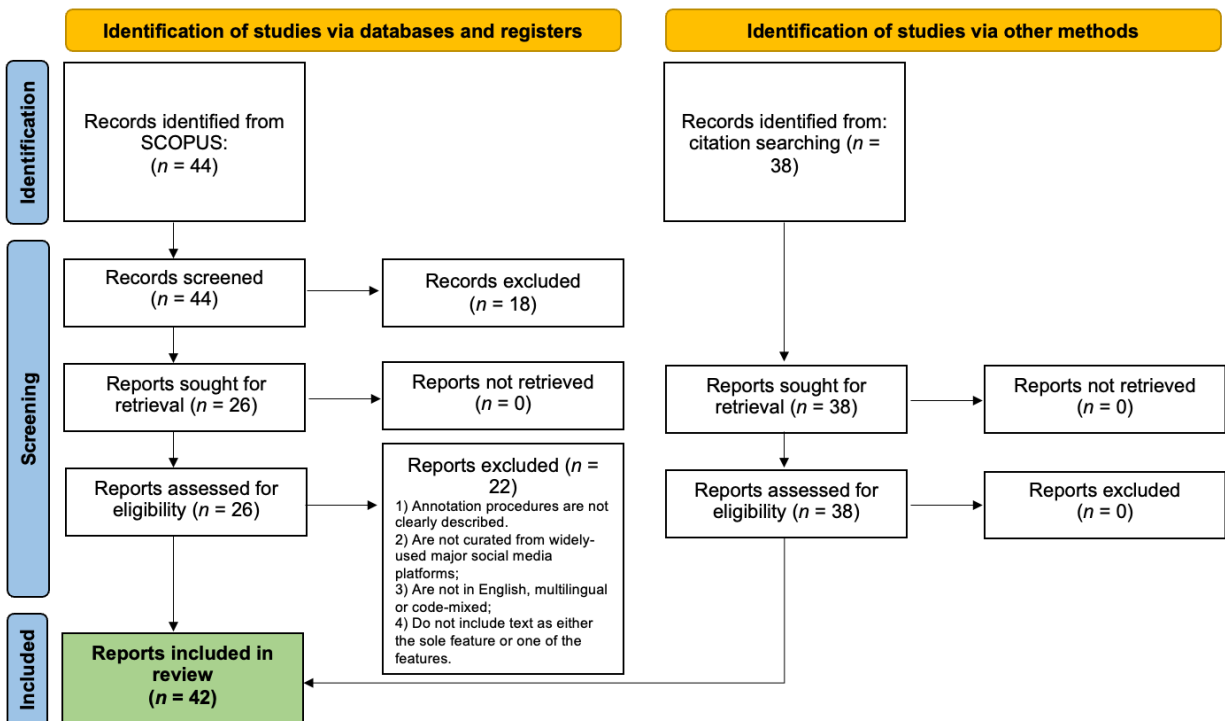5) The dataset includes textual content.



*Figure C1. PRISMA diagram for the selection of papers presenting labeled datasets for online abuse.*

---

[1] Focusing on monolingual settings allows us to address these issues directly before extending analyses to multiple languages, where cultural variations and linguistic nuances further complicate intent inference.

*Query for algorithm papers*

KEY ("social media" AND ("NLP" OR "Natural Language Processing") AND "de-taction" AND ("hate speech" OR "abus*" OR "offens*" OR "cyberbully*")) OR TITLE ("social media" AND "detection" AND ("hate speech" OR "abus*" OR "offens*" OR "cyberbully*")) AND (LIMIT-TO (LANGUAGE, "English"))

*Inclusion criteria for algorithm papers*

Similar to our survey of datasets, we removed papers that were not accessible, not written in English, or which did not describe an algorithm for detection of online abuse. We removed models designed for multilingual or non-English tasks.
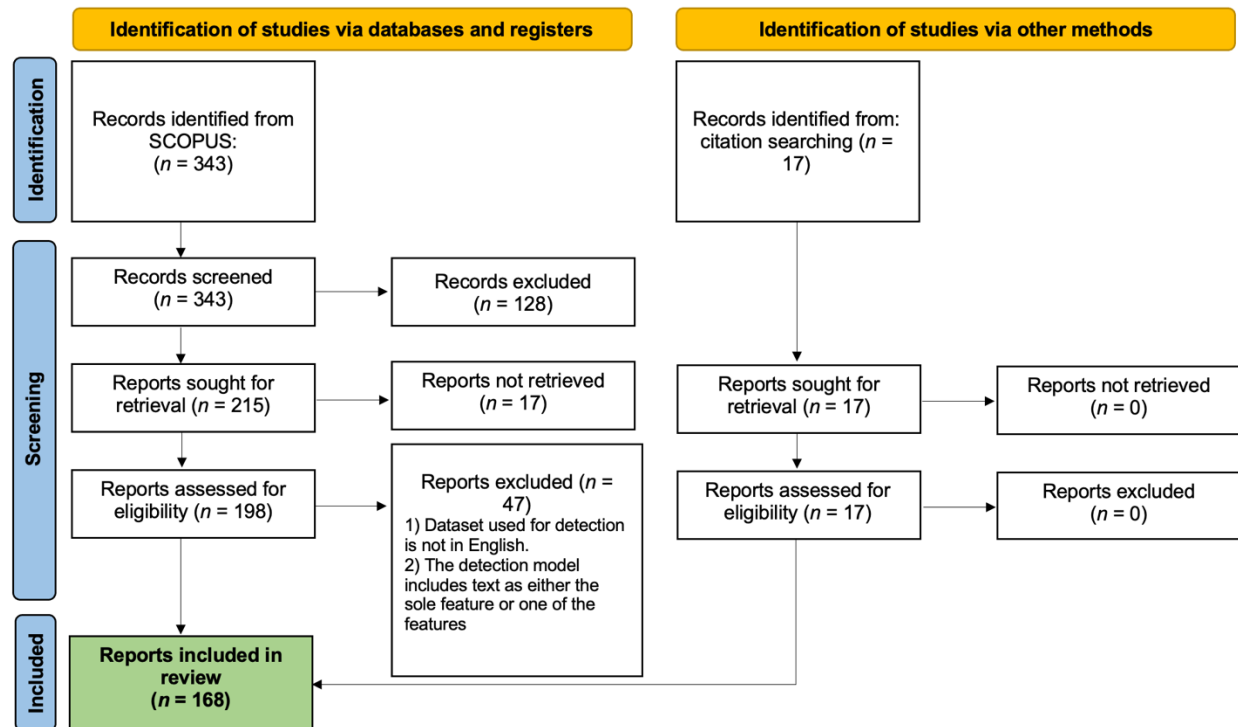


*Figure C2. PRISMA diagram for the selection of papers presenting detection algorithms for online abuse.*