

AggregHate: An Efficient Aggregative Approach for the Detection of Hatemongers on Social Platforms

Anonymous ACL submission

Abstract

The automatic detection of online hate speech serves as a crucial step in the detoxification of the online discourse. Moreover, accurate classification can promote a better understanding of the proliferation of hate as a social phenomenon. While most prior work focus on the detection of hateful *utterances*, we argue that focusing on the *user* level is as important, albeit challenging. In this paper we consider a multimodal aggregative approach for the detection of hate-mongers, taking into account the potentially hateful texts, user activity, and the user network. We evaluate our methods on three unique datasets X (Twitter), Gab, and Parler showing that a processing a user’s texts in her social context significantly improves the detection of hate mongers, compared to previously used text and graph-based methods. Our method can be then used to improve the classification of coded messages, dog-whistling, and racial gas-lighting, as well as inform intervention measures. Moreover, our approach is highly efficient even for very large datasets and networks.

Disclaimer: The illustrative examples in Table 1 may be offensive to some readers.

1 Introduction

The rising popularity of social platforms enhanced the hateful content targeting minorities and led to the proliferation of online hate speech (Waseem and Hovy, 2016; Laub, 2019). Accordingly, there is a growing body of research on the appearance and magnitude of hate speech on social media, in general, (Knuttila, 2011; Chandrasekharan et al., 2017; Zannettou et al., 2018), and on hate speech detection, in particular, (Saleem et al., 2017; Waseem and Hovy, 2016; Davidson et al., 2017). Hate

speech is not merely an online inconvenience as shooting, bombing, stabbing, beating, and vandalism are often linked to online activity (Munn, 2019; Malevich and Robertso, 2019; Thomas, 2019; McIlroy-Young and Anderson, 2019; Mathew et al., 2019; ADL, 2023).

Hate is not promoted by isolated individuals, but rather by communities that often exist within larger communities. Shifting the focus from the post level to the user level can be beneficial on various levels (pun intended): it provides a better understanding of the dynamics of hate groups and the ways hate infects individuals and communities; it could be reapplied in the post level, improving the detection of coded language, dog whistling and gas-lighting; finally, it could provide the social structure in which the impact of different intervention techniques¹ could be tested (Thomas and Wahedi, 2023).

Many of the challenges inherent to hate speech classification on the utterance level are discussed by Talat et al. (2018); Fortuna et al. (2022); Velankar et al. (2022); Kim et al. (2024), among others. A recent survey of the applicability of state-of-the-art LLMs highlights the shortcomings of the use of LLMs for nuanced tasks such as hate, toxicity and empathy (Ziems et al., 2024).

We provide some illustrative examples in Table 1, complemented with a short analysis, demonstrating the specific challenges and outlining the possible benefits of shifting our attention from the utterance (textual) level to the speaker (user) level.

The explicit racism of utterance #1 is detected by all models, however, deHateBERT

¹Intervention alternatives depend on legal regulation and platform policy and range from the removal of a specific post, the flagging of a specific post, to the temporal or indefinite suspension of an account.

| | Text | DHB | JS | FTDB | GPT4 | Gemini 1.5Pro |
|---|---|-----|-----|------|------|---------------|
| 1 | "Muslim terrorists did this and obama brought them here to run our government. Gotta get rid of Obammy, Omar and all these other fucking Muslims" | .57 | .9 | .98 | ✓ | ✓ |
| 2 | "is there a more debased creature in our modern pop-cultural cesspool than (((Amy Schumer)))? Probably not." | .22 | .35 | .92 | ✓? | ✓ |
| 3 | "We must secure the existence of our people and a future for white children!" | .03 | .36 | .02 | ✓ | ✓? |
| 4 | "If white people are going to continue existing - we must support our race" | .06 | .4 | .69 | ✓? | ✓? |
| 5 | "When the youngest white generation becomes adult they will be 5% of the world population, down from 1/3 in 1900" | .06 | .25 | .01 | ✗ | ✓? |
| 6 | "Did you know that most dog breeds were bred by whites? We're awesome and they're awesome!" | .2 | .25 | .65 | ✗ | ✗ |
| 7 | "White lives matter!" | .02 | .24 | .02 | ✗ | ✗? |
| 8 | "Blue lives matter!" | .02 | .01 | .01 | ✗ | ✗? |

Table 1: Examples of hate-promoting texts. DHB: the hate-likelihood of the text assigned by deHateBERT (Aluru et al., 2020); JS: hate prediction by Google’s Jigsaw; FTDB: DistillBERT (Sanh et al., 2019a) fine-tuned on our datasets; GPT-4 and Gemini 1.5 Pro, using the prompt “Which of the following texts should be classified as hate speech?”. (“Which of the following utterances may be perceived as hate speech?”). A ‘?’ postscript indicates that a nuanced prediction was generated by the model (see full generated predictions in Appendix A.

(Aluru et al., 2020) (DHB) – a BERT model fine-tuned for the detection of hate speech – assigned a likelihood of only 0.57 for it to be classified as hate speech. Three of the models, ChatGPT-4 included, failed to detect the antisemitism conveyed in utterance #2 as it requires “external” knowledge, namely that Schumer is of Jewish heritage and that the (((echo))) symbol is a (newly) recognized hate symbol². The text in utterance #3, known as the ‘14 words’ – “the most popular white supremacist slogan in the world”³. However, it may appear innocent to the uninformed (human) eye, as well as to models that did not encounter it (in context) in training. As a result, it was identified only by ChatGPT-4 and Gemini 1.5 Pro. Posts #4-7 are thinly veiled allusions to the “14 words”. The fine-tuned distilBERT (FTDB) captures #4 and #6 but misses the “14 words”. Both #7 and #8 twist the slogan ‘Black Lives Matter’ – originally a protest against police brutality. While #7 carries an explicit white-supremacist tone⁴, #8⁵ is not recognized as hate-speech, although it is often used alongside racial slurs. Even state-of-the-art models such as ChatGPT-4 and Gemini 1.5

²www.adl.org/resources/hate-symbol/echo

³www.adl.org/resources/hate-symbol/14-words

⁴www.adl.org/resources/hate-symbol/white-lives-matter

⁵A response to the BLM movement, suggesting that attacking policemen should be considered a hate crime.

Pro either completely fail to detect the hate conveyed in utterances 4-8, or flag it with a low confidence.

While all of the texts in Table 1 pose a challenge to computational models and humans alike, contextual information and careful aggregation can be used to achieve classification on the user level: multiple implicit posts (coded, ambiguous, dog-whistling, or gas-lighting) posted by a single user can reinforce the weak signal obtained from a single post. Similarly, association with other users or explicit posts can reinforce a weak signal of a context-less post.

Contextual aggregation, however, is not straightforward as it depends on a number of philosophical and practical considerations stemming from the likelihood or the confidence of a model (or a human) to assign the class for a specific post. For example, one can argue that a single, though explicit, hateful post may not be enough to label the user as a racist or a hate-monger. However, even if a single explicit post is sufficient – what would be the user label in case of two *implicit* posts, or a few dozen of posts, each is predicted to be hateful with a low confidence? Looking at the illustrative posts in Table 1, we propose a principled way to effectively combine predictions and modalities in order to achieve an accurate classification.

To this end we explore three fundamentally

different approaches for contextual aggregation: (i) using binary weights with a fixed threshold, (ii) using a relational aggregation conditioned on the social context, and (iii) using sidtributional aggregation conditioned on aggregated confidence levels. Finally, we combine these methods to form a multimodal classification model.

Contribution Our contribution of this work is threefold:

1. We propose a robust and efficient multimodal aggregative approach for the detection of hate-mongers.
2. We demonstrate the benefits of contextual aggregation over three unique datasets (Twitter, Gab, and Parler).
3. We share a novel annotated dataset of Parler hate.

2 Related Work

A comprehensive overview of methods and benchmarks for hate-speech detection is provided by [Alkomah and Ma \(2022\)](#), while many of the challenges in current approaches are surveyed by [ElSherief et al. \(2021\)](#); [Velankar et al. \(2022\)](#); [Fortuna et al. \(2022\)](#), among others. Subjectivity and the incomplete definition of hate-speech are addressed by [Khurana et al. \(2022\)](#) while the limitations of transfer learning for the task are demonstrated by [Israeli and Tsur \(2022\)](#), and a set of functional tests to evaluate the performance of different models trained on different benchmarks, types, targets, and languages were proposed by [Röttger et al. \(2021, 2022\)](#).

A taxonomy of implicit hate was developed and shared by [ElSherief et al. \(2021\)](#). Unfortunately, most of the annotated tweets are no longer available, thus user level aggregation cannot be achieved.

A growing number of works has shifted the attention from the utterance level to the user level. [Waseem and Hovy \(2016\)](#) analyze the relation between demographic features and hate speech, while [Ribeiro et al. \(2017\)](#) explores the differences between account meta-features of hateful and non-hateful Twitter users. Both works are of an exploratory nature.

[Arviv et al. \(2021\)](#) detects hate mongers using a multi-modal architecture that combines

three streams of post-level predictions: the tweets of the target user, her followers, and her followees.

A two-step approach considering both the textual and the network modalities was proposed by [Ribeiro et al. \(2018\)](#) and extended by [Israeli and Tsur \(2022\)](#): In the first step seed nodes (users) are detected based on the textual signal (keyword matching in [Ribeiro et al.](#) and a fine-tuned BERT in [Israeli and Tsur](#)). In the second step a diffusion model is applied in order to propagate the initial hate assignments across the social network.

A number of works use Graph Neural Networks (GNN) in order to detect hate speech or hate mongers. [Li et al. \(2021\)](#) presented HateGNN, using textual similarity and the appearance of predefined hate terms as part of the objective that produces the node embeddings to be used for classification. The learned embeddings depend heavily on the training data, the type of hate, and the supervision. [Miao et al. \(2022\)](#) proposed an end-to-end framework, enriching a BERT classifier with Graph Attention Networks. However, in spite of the use of graph networks, this approach is applied for the detection of hate only on the post level rather than on the user level. Several supervised, unsupervised and semi supervised models were explored by ([Das et al., 2021](#)), including the state-of-the-art GNN models such as AGNN ([Thekumparampil et al., 2018](#)) in order to classify users as hateful and non-hateful.

[Nirmal et al. \(2024\)](#) proposed *SHIELD* - a framework that leverages LLM-extracted rationales to augment a base hate speech detection model to facilitate faithful interpretability. However, it was demonstrated that fine-tuned classification models outperform state-of-the-art LLMs over tasks that involve social nuances, e.g., detection of humor, empathy, toxicity and hate ([Ziems et al., 2024](#)).

3 Multimodal Aggregative Approaches

3.1 Aggregative Approaches

Utterance-level Classification (C^T) The basic building block of all aggregative approaches is the classification of a single utterance. Any classification model can be used for the utterance level as long as its output

can be interpreted as a probability (e.g., by applying the sigmoid function). The probability is needed in order to assign a “confidence” threshold, before making the binary decision. Formally:

$$C^T(t) = \begin{cases} 1 & \theta(t) \geq \tau^T \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Where t is a text snippet (post, tweet, utterance), θ is a classification model, and the hyper parameter τ^T is the post-level threshold reflecting the sensitivity of the model (or the community) to implicit forms of hate speech. In this work we follow [Israeli and Tsur \(2022\)](#) and use a DistilBERT classifier ([Sanh et al., 2019b](#)) fine-tuned for hate-speech detection on our three datasets (see Section 4).

User-level Classification (C^U) The detection of hateful *users* is inherently related to the user’s posts. Given a user u and a the user level threshold τ^U , the generic user classification function is given by:

$$C^U(u) = \begin{cases} 1 & \Theta(u) \geq \tau^U \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Naive Aggregation with Fixed Threshold

Given T^u – the set of posts published by u , a naive aggregative approach could simply count the number of hateful posts (decided by the C^T) published by a user u :

$$\Theta(u) = \sum_{t \in T^u} C^T(t) \quad (3)$$

While naive aggregation does not require any training, the threshold τ^U can be used to control the sensitivity of the model: setting $\tau^U = 1$ implies zero-tolerance (often too harsh⁶) approach, while setting $\tau^U = a$, $a \in \{10, 20, 30, \dots\}$ implies a more conservative approach, requiring a user to be a “repeated offender” in order to be labeled as a hate-monger. For convenience, we denote to the naive aggregation (fixed threshold) Θ_F , and the naive classification function C_F .

In order to achieve flexibility and robustness one should consider other aggregation approaches that account for the nuanced and

⁶Remember that in the standard case $\tau^T = 0.5$ thus a user will be labeled a hate-monger even if $\theta(t) = 0.51$ for one of her posts and $\theta(t') < 0.5 \quad \forall t' \in T^u / \{t\}$.

implicit utterances as well as for the social context in which a user is embedded. The different aggregative approaches could be combined to a user feature vector and Θ could be trained to optimize the weight of each feature. Specifically, instead of a naive aggregation we consider two types of aggregative functions: Relational Aggregation (Θ_R) and Distributional Aggregation (Θ_D). Given an annotated dataset, we train Θ to optimize C^U , where Θ is a logistic regression classifier.

In the remainder of this section we motivate the use of each aggregative approach and provide the formal definition of the model.

Relational Aggregation The Aristotelian concept of man as a social and political animal informed decades of sociological research concerning the importance of community to individual identity ([McMillan and Chavis, 1986](#); [Wellman and Gulia, 1999](#)). Unfortunately, identity is often forged by association with hate groups, off and on-line ([Gordon, 2017](#); [Govers et al., 2023](#)). It is therefore reasonable to assume that hateful content circulating in u ’s ego network should be taken into account and inform the label assigned to u . That is, even if a conservative τ^U is used, the associates of u may push her over the threshold.

Formally, given $G(V, E)$ – a directed social network where V is the set of users and $(u, v) \in E$ indicates a directed edge $u \rightarrow v$, we define \overleftarrow{u} and \overrightarrow{u} as the sets of followers and followees of u , respectively. In the relational case, $\Theta_R(u)$ is a linear combination of three terms – the naive aggregation over u ’s texts and the percentages of hateful users among his followers and followees:

$$\Theta_R(u) = \alpha \cdot C_F(u) + \beta \cdot \frac{1}{|\overleftarrow{u}|} \sum_{v \in \overleftarrow{u}} C_F(v) + \gamma \cdot \frac{1}{|\overrightarrow{u}|} \sum_{v \in \overrightarrow{u}} C_F(v) \quad (4)$$

The values of α , β , and γ are optimized through training, essentially reflecting the importance of the user’s posts and the posts of his followers and followees.

Distributional Aggregation While relational aggregation takes the social context into account, it does not address variations the intensity of the promoted hate in terms

of “comitement” (what percentage of a user’s stream is hate) or implicitness (e.g., the user tries to tread the fine line without violating the platform rules). For example, consider two users u and v posting utterances #3–#8 in Table 1. However, imagine that these are the only posts published by u , while v published hundreds of other posts, non of which is hateful. Intuitively, one may argue that v ’s questionable posts are diluted, compared to u ’s, thus we should require a higher level of confidence in order to assign him the ‘hate monger’ label (and compare to a user v ’ posting the more explicit #1 & #2 along with hundreds of non-hateful posts).

In order to address these variations we use distributional aggregation: instead of counting the user’s hateful posts (Θ_F), we look at the distribution of the user’s posts as a k -dimensional vector where k determines the number of bins used to approximate the probability density function.

Specifically, we consider two distributional spaces – bins and quantiles. In the bin-based representation the $[0,1]$ range is divided to k equal bins, each bin (entry in the vector) holds the number of posts with the corresponding hate score assigned by $\theta(t)$. In the quantile-based representation the k bins are unique for each user u as they are defined over the range $[\min(\theta(t)), \max(\theta(t))]$ for $t \in T^u$. We therefore want to optimize

$$\Theta_D(u) = \sum_{i=1}^k w_i \cdot \sigma(B_i(u)) \quad (4)$$

where $B_i(u)$ denotes the number of utterances for which $\theta(t)$ falls in the i^{th} bin and σ is the softmax function.

Multimodal Aggregation Finally, the different aggregation methods can be used together in a multimodal manner – combining the hate-score distribution of the user utterances and the hate levels in the user’s ego network. Using Θ_D^b and Θ_D^q to denote the bin-based and quantile-based representations we can formulate the combined model as

$$\Theta(u) = \Theta_D^b(u) + \Theta_D^q + \Theta^R$$

3.2 Social-aware Baselines

In this section we briefly describe five strong baseline algorithms we use for comparison.

These five algorithms, all leverage the structure of the social network and have proved useful in an array of node classification tasks, including the detection of hate speech and hateful users.

DeGroot’s Diffusion The DeGroot’s model is a simple yet strong framework to classify nodes in a network through belief propagation. The model achieved good results in the detection of hateful users on Gab and Parler (Ribeiro et al., 2018; Israeli and Tsur, 2022).

Graph Neural Networks (GNNs) Das et al. (2021) explored several GNN methods for detecting hateful users on Gab and Twitter. The network embeddings incorporate nodal features, e.g., textual representations of the user’s utterances in learning the node and graph embeddings. In this work we use the following four algorithms:

1. **GCN:** Graph Convolutional Networks (Kipf and Welling, 2016) uses localized approximation in learning node embeddings through convolutional layers.
2. **GAT:** Graph Attention Network learn nodal representations by combining features of the nodes in the ego network of the focal node, setting their importance via attention layers.
3. **GraphSAGE:** The Graph Sample and Aggregate (Hamilton et al., 2018) learns nodal embeddings by sampling a predefined number of nodes from the ego network of a focal node u , then applying attention layers in order to assign the importance of each node in the aggregative representation.
4. **AGNN:** The Attention-based Graph Neural Network (Thekumparampil et al., 2018) uses attention layers to learn a dynamic and adaptive local summary of the neighborhood of each focal node.

4 Datasets and Annotation

We evaluate our aggregative approaches over three very different datasets: Twitter-Echo, Gab, and Parler. The number of posts, users, and label breakdown for each dataset are provided in Table 2. The remainder of this section provides further details regarding the datasets and our annotation process.

| Dataset | Source | Raw Data | | Annotated Data | | | |
|----------------|---------------------------------------|-------------------|--------------------|-------------------|--------------------|--------|--------|
| | | #Posts | #Users | #Posts | % Hate | #Users | % Hate |
| Echo (Twitter) | Arviv et al. (2021) | 18M | 7.07K | 4630 | 8.2% | 1000 | 15.4% |
| Gab | Arviv and Tsur (2021) | 22M | 336.7K | 5000 | 5.1% | 1000 | 24.8% |
| Parler | This work | 183M [†] | 4.08M [†] | 8262 [‡] | 32.9% [‡] | 890 | 25.4% |

Table 2: Datasets statistics. The raw Parler data (marked †) were shared by ([Aliapoulios et al., 2021](#)) and the post-level annotations (‡) were shared by ([Israeli and Tsur, 2022](#)). User-level annotations are shared as part of this paper.

Echo (Twitter) The triple parentheses, or triple brackets, also known as the (((echo))), is an antisemitic symbol that is used to highlight the names of individuals of Jewish background (e.g., actress and comedian Amy Schumer, see utterance #2 in Table 1), organizations owned by Jewish people (e.g., Ben & Jerry’s), or organizations accused of promoting “Jewish globalist values” (e.g., the International Monetary Fund). The Echo dataset curated by [Arviv et al. \(2021\)](#) contains over 18M English tweets posted by ~7K echo users between May and June 2016. Annotations are provided at the tweet and the user level. An important feature of this dataset is that all users have utterances containing the echo symbol, although some users use it in a non-hateful manner, e.g., to symbolize a hug. This ambiguous nature of the symbol makes hate detection challenging.

Gab Gab, launched on August 2016, was created as an alternative to Twitter, positioning itself as putting “people and free speech first”, welcoming users suspended from other social networks. Gab posts (called *gabs*) are limited to 300 characters, and users can repost, quote or reply to previously created gabs. Gab permits pornographic and obscene content, as long as it is labeled *NSFW* (‘not safe for work’).

The raw Gab dataset was introduced by [Zannettou et al. \(2018\)](#). It was collected using Gab’s API with the snowball methodology. More specifically, the researchers obtained data for the most popular users as returned by Gab’s API and iteratively collected data from all of their followers and their followees. They collected three types of information: basic details about Gab accounts (including username, score, and date of account creation); all the posts for each Gab user in the dataset; and all the followers and followees

of each user, which allow the reconstruction of a social network. Overall, this dataset contains 22.1M posts from 336.8K users, posted between August 2016 and January 2018. [Arviv and Tsur \(2021\)](#) shared an annotated dataset based on the aforementioned above, containing 60K labeled posts and 1K labeled users.

Parler Alluding to the French verb ‘to speak’, Parler was launched on August 2018.⁷ The platform branded itself as “The World’s Town Square” a place to “*Speak freely and express yourself openly, without fear of being “deplatformed” for your views*”⁸.

Parler users post texts (called *parlays*) of up to 1000 characters. Users can reply to parlays and to previous replies. Parler supports a reposting mechanism similar to Twitter’s retweets (called ‘echos’, not to confuse with the (((echo))) hate symbol, see above). Parler’s official guidelines⁹ explicitly allowed “trolling” and “not-safe-for-work” (NSFW) content, include only three “principles” prohibiting “unlawful acts”, citing “Obvious examples include: child sexual abuse material, content posted by or on behalf of terrorist organizations, intellectual property theft”.

[Aliapoulios et al. \(2021\)](#) presented a dataset of tens of millions of Parler messages. [Israeli and Tsur \(2022\)](#) used this dataset to introduce an annotated dataset for hate speech (post level). Their 10K dataset consists of 3224 posts (32.8%) labeled as hateful and 6053 (59.8%)

⁷On April 2023 the platform was acquired by Starboard and was taken offline to “undergo a strategic assessment” (Starboard announcement on Parler’s landing page <https://parler.com/>, accessed: 5/8/2023). The platform was relaunched in February 2024 announcing it is “breaking free from the constraints of conventional platforms” (accessed: 6/5/2024).

⁸Parler branding on its landing page (accessed: 3/10/2022)

⁹<https://parler.com/documents/guidelines.pdf> (accessed: 4/17/2022)

496 as non-hateful. However, a Parler *user-based*
497 annotated dataset has yet to be introduced.
498 Hence, as part of this research, we create the
499 first annotated dataset of Parler users.

500 **Annotation of Parler Users** The Parler
501 dataset presented by Aliapoulios et al. (2021)
502 consists of $\sim 4M$ users. We focus on a subset of
503 users matching the following criteria: (i) The
504 account exists for at least six months; (ii) The
505 user showed some activity (posted at least 30
506 posts); (iii) The primary language of the user is
507 English. These three rules left us with a subset
508 of users, denoted U^* , from which we sampled
509 users for annotation. Following the protocol
510 used by Ribeiro et al. (2018) and Israeli and
511 Tsur (2022), we used stratified subsampling
512 mitigate bias (most users and vast majority of
513 posts are not hateful).

514 Annotation was done by 94 senior year
515 Data Science B.Sc students for bonus course
516 credit. Annotators were introduced to Parler
517 and were given explicit instructions about
518 the annotation task. The annotation process
519 involved rating each account on a 1–5 scale
520 (non-hateful – extremely hateful). We ensured
521 that each user is annotated by three annota-
522 tors. The full annotation guidelines and fur-
523 ther details regarding the annotation protocol
524 are available in Appendix C.

525 5 Results and Analysis

526 5.1 Experimental Settings

527 For the utterance-level classification, we fine-
528 tune DistilBERT (Sanh et al., 2019a) on each
529 datasets. We used a batch size of 32, a maxi-
530 mum number of epochs of 20 and a validation
531 split of 0.2. We also used an early stopping
532 callback with a patience of 5 epochs on the
533 validation loss.

534 For the user-level classification, we use 5-
535 Fold cross-validation for all of the methods.
536 For each dataset, we considered the largest
537 (weakly) connected component. Specifically,
538 for the GNN methods we used the same ex-
539 perimental settings as proposed by Das et al.
540 (2021). Appendix B provides further details
541 regarding the full networks, the number of
542 connected components and the statistics re-
543 garding the largest connected component in
544 each dataset.

545 For the DeGroot’s method, we followed the
546 protocol used by Israeli and Tsur (2022).

547 5.2 Results

548 **Utterance Level Prediction** All aggregative
549 models require an initial step of classification
550 of individual utterances. While the focus of
551 this paper is the classification on the *user* level,
552 we report the results on the utterance level in
553 order to highlight the challenge and the benefit
554 of the aggregative methods.

555 Results for each dataset are presented in Ta-
556 ble 4. Note the very low F-score achieved ob-
557 the Gab dataset and compare to the significant
558 improvement achieved by shifting to the user
559 level (below). The density of the mean hate
560 score per user is presented in Figure 1, high-
561 lighting the differences between the datasets.

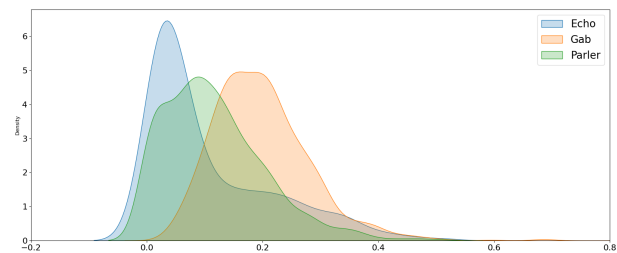


Figure 1: The Density of the mean utterance scores per user for each Dataset

562 **User level results** Detailed results of each
563 method over each of the three datasets are
564 presented in Table 3.

565 Looking at the F-score achieved by the dif-
566 ferent methods we observe that our aggrega-
567 tive approach consistently outperform the
568 baseline algorithms. While the multimodal
569 aggregation yields the best results in the Echo
570 and Gab datasets, it is ranked third on the
571 Parler dataset, with the relational aggregation
572 ranked first. We speculate that this is may be
573 attributed to the size and the unique charac-
574 teristics of the Parler network (see Appendix B).
575 A thorough analysis of the social networks is
576 out of the scope of this paper.

577 Breaking the multimodal aggregative model
578 to separate aggregative models (Relational,
579 Distributional-bins, Distributional-quantiles
580 and bins and quantiles combined) shows that
581 in most cases all these approaches are com-
582 petitive with the top performing models and
583 outperform the strong baselines.

| | Method | Precision | Recall | F1 | ROC AUC |
|--------|--|----------------------|----------------------|----------------------|----------------------|
| Echo | DeGroot’s Diffusion | 0.472 ± 0.389 | 0.255 ± 0.261 | 0.320 ± 0.310 | 0.610 ± 0.122 |
| | GCN | 0.443 ± 0.118 | 0.914 ± 0.096 | 0.585 ± 0.104 | 0.797 ± 0.102 |
| | GraphSAGE | 0.629 ± 0.124 | 0.950 ± 0.041 | 0.752 ± 0.088 | 0.944 ± 0.024 |
| | GAT | 0.548 ± 0.131 | 0.653 ± 0.241 | 0.574 ± 0.123 | 0.772 ± 0.091 |
| | AGNN | 0.759 ± 0.072 | 0.914 ± 0.054 | 0.826 ± 0.035 | 0.963 ± 0.032 |
| | Fixed-Threshold | 0.654 ± 0.063 | 0.627 ± 0.095 | 0.633 ± 0.040 | 0.836 ± 0.040 |
| | Relational Aggregation | 0.820 ± 0.060 | 0.834 ± 0.055 | 0.825 ± 0.042 | 0.956 ± 0.013 |
| | Distributional (bins) | 0.772 ± 0.045 | 0.871 ± 0.058 | 0.817 ± 0.042 | 0.944 ± 0.019 |
| | Distributional (quantiles) | 0.747 ± 0.064 | 0.899 ± 0.047 | 0.815 ± 0.053 | 0.942 ± 0.018 |
| | Distributional (bins+quantiles) | 0.757 ± 0.058 | 0.885 ± 0.052 | 0.815 ± 0.049 | 0.946 ± 0.021 |
| | Multimodal (relational+bins+quantiles) | 0.781 ± 0.023 | 0.899 ± 0.048 | 0.836 ± 0.028 | 0.961 ± 0.011 |
| Gab | DeGroot’s Diffusion | 0.314 ± 0.001 | 0.777 ± 0.000 | 0.447 ± 0.001 | 0.604 ± 0.000 |
| | GCN | 0.241 ± 0.109 | 0.678 ± 0.428 | 0.334 ± 0.169 | 0.594 ± 0.059 |
| | GraphSAGE | 0.317 ± 0.066 | 0.559 ± 0.206 | 0.388 ± 0.077 | 0.582 ± 0.077 |
| | GAT | 0.194 ± 0.115 | 0.457 ± 0.387 | 0.264 ± 0.174 | 0.501 ± 0.060 |
| | AGNN | 0.340 ± 0.031 | 0.600 ± 0.190 | 0.423 ± 0.030 | 0.679 ± 0.024 |
| | Fixed-Threshold | 0.497 ± 0.078 | 0.351 ± 0.066 | 0.411 ± 0.070 | 0.722 ± 0.039 |
| | Relational Aggregation | 0.408 ± 0.061 | 0.437 ± 0.086 | 0.419 ± 0.063 | 0.675 ± 0.043 |
| | Distributional (bins) | 0.461 ± 0.034 | 0.649 ± 0.044 | 0.538 ± 0.024 | 0.763 ± 0.016 |
| | Distributional (quantiles) | 0.429 ± 0.027 | 0.702 ± 0.056 | 0.532 ± 0.033 | 0.770 ± 0.018 |
| | Distributional (bins+quantiles) | 0.435 ± 0.026 | 0.714 ± 0.043 | 0.540 ± 0.029 | 0.769 ± 0.016 |
| | Multimodal (relational+bins+quantiles) | 0.452 ± 0.027 | 0.702 ± 0.042 | 0.550 ± 0.032 | 0.777 ± 0.020 |
| Parler | DeGroot’s Diffusion | 0.395 ± 0.221 | 0.441 ± 0.247 | 0.417 ± 0.233 | 0.644 ± 0.081 |
| | GCN | 0.284 ± 0.054 | 0.760 ± 0.404 | 0.348 ± 0.130 | 0.644 ± 0.145 |
| | GraphSAGE | 0.309 ± 0.092 | 0.649 ± 0.189 | 0.394 ± 0.028 | 0.497 ± 0.061 |
| | GAT | 0.379 ± 0.051 | 0.731 ± 0.164 | 0.488 ± 0.013 | 0.746 ± 0.044 |
| | AGNN | 0.369 ± 0.081 | 0.552 ± 0.255 | 0.416 ± 0.082 | 0.667 ± 0.069 |
| | Fixed-Threshold | 0.470 ± 0.050 | 0.369 ± 0.040 | 0.412 ± 0.035 | 0.693 ± 0.026 |
| | Relational Aggregation | 0.519 ± 0.073 | 0.509 ± 0.082 | 0.513 ± 0.074 | 0.730 ± 0.051 |
| | Distributional (bins) | 0.284 ± 0.037 | 0.500 ± 0.075 | 0.362 ± 0.049 | 0.575 ± 0.024 |
| | Distributional (quantiles) | 0.324 ± 0.018 | 0.734 ± 0.057 | 0.449 ± 0.025 | 0.611 ± 0.040 |
| | Distributional (bins+quantiles) | 0.324 ± 0.021 | 0.738 ± 0.051 | 0.450 ± 0.027 | 0.618 ± 0.034 |
| | Multimodal (relational+bins+quantiles) | 0.370 ± 0.023 | 0.680 ± 0.082 | 0.478 ± 0.035 | 0.699 ± 0.044 |

Table 3: 5-Fold CV results on the test sets of Echo, Gab, and Parler datasets using the best (F1-score-wise) configuration.

| Dataset | Precision | Recall | F1 Score |
|---------|-----------|--------|----------|
| Echo | 0.4122 | 0.8026 | 0.5446 |
| Gab | 0.2057 | 0.5472 | 0.2990 |
| Parler | 0.6316 | 0.8177 | 0.7127 |

Table 4: Performance metrics of the utterance-level model for Echo, Gab, and Parler Datasets

584 Interestingly, looking at the results of the
585 relational method, we observe that the im-
586 portance of the different components differ
587 across datasets: best F1-score was achieved
588 using $\alpha = 0.608$, $\beta = 0.776$, $\gamma = 1.467$ for the
589 Echo dataset; $\alpha = 0.776$, $\beta = 0.085$, $\gamma = 0.108$
590 for the Gab dataset and $\alpha = 0.239$, $\beta = 0.254$,
591 $\gamma = 0.24$ for Parler. This result highlights the
592 importance of the network structure and the
593 dynamics and norms of each platform. This is
594 also evident from the distribution of the mean
595 utterance score in each platform.

6 Conclusion

596 We proposed a robust and efficient multi-
597 modal aggregative method, combining text
598 and social context through relational and dis-
599 tributional aggregations. We demonstrated
600 the benefits of this approach for the task of
601 hate speech and hatermonger detection over
602 three unique and very different datasets from
603 three social platforms: X (Twitter), Gab and
604 Parler.
605

606 Future work takes three trajectories: (i) Im-
607 proving the integration of the different modal-
608 ities, and (ii) Better understanding the unique
609 characteristics of each of the platforms and
610 the datasets and the ways they impact the re-
611 sults, and (iii) Reapplication of the aggregative
612 methods to improve classification of evasive
613 utterances and emerging forms of coded hate
614 speech.

| | | | |
|-----|--|--|-----|
| 615 | 7 Limitations and Ethical | | |
| 616 | Considerations | | |
| 617 | 7.1 Limitations | | |
| 618 | This work has a number of limitations, two of | | |
| 619 | which are listed below (and other are referred | | |
| 620 | to in the directions for future work): | | |
| 621 | (i) All three aggregation procedures depend | | |
| 622 | on the basic classifier $\theta(t)$. An evasive user | | |
| 623 | that is careful with his words, using only | | |
| 624 | coded language or consistent gas-lighting may | | |
| 625 | not be identified at all if $\theta(t) < \tau^T$ for all | | |
| 626 | $t \in T^u$ for a user u . | | |
| 627 | (ii) We use the fixed threshold τ^u as a | | |
| 628 | strong baseline, however, a more comprehen- | | |
| 629 | sive comparison to other existing methods | | |
| 630 | (e.g., diffusion-based and GNNs) should be | | |
| 631 | considered in future work. | | |
| 632 | 7.2 Ethical Considerations | | |
| 633 | The detection of hate mongers on social plat- | | |
| 634 | forms and the detoxification of the discourse | | |
| 635 | are of great importance. However, wrongly | | |
| 636 | labeling a user, may have severe emotional | | |
| 637 | and legal implications. In this work, we ex- | | |
| 638 | plorate computational approaches that can be | | |
| 639 | used to study the phenomena at scale. As a | | |
| 640 | practical policy, user labeling should be done | | |
| 641 | carefully, ideally involving a ‘man-in-the-loop’ | | |
| 642 | procedure. | | |
| 643 | References | | |
| 644 | Center on Extremism ADL. 2023. Shooter in allen, | | |
| 645 | texas, embraced antisemitism, misogyny and | | |
| 646 | white supremacy. | | |
| 647 | Max Aliapoulios, Emmi Bevensee, Jeremy Black- | | |
| 648 | burn, Barry Bradlyn, Emiliano De Cristofaro, Gi- | | |
| 649 | anluca Stringhini, and Savvas Zannettou. 2021. | | |
| 650 | An early look at the parler online social network. | | |
| 651 | <i>arXiv preprint arXiv:2101.03820.</i> | | |
| 652 | Fatimah Alkomah and Xiaogang Ma. 2022. A lit- | | |
| 653 | erature review of textual hate speech detection | | |
| 654 | methods and datasets. <i>Information</i> , 13:273. | | |
| 655 | Sai Saketh Aluru, Binny Mathew, Punyajoy Saha, | | |
| 656 | and Animesh Mukherjee. 2020. Deep learning | | |
| 657 | models for multilingual hate speech detection. | | |
| 658 | <i>arXiv preprint arXiv:2004.06465.</i> | | |
| 659 | Eyal Arviv, Simo Hanouna, and Oren Tsur. 2021. | | |
| 660 | It’s a thin line between love and hate: Using the | | |
| 661 | echo in modeling dynamics of racist online com- | | |
| 662 | munities. In <i>Proceedings of the International AAI</i> | | |
| | <i>Conference on Web and Social Media</i> , volume 15, | | 663 |
| | pages 61–70. | | 664 |
| | Eyal Arviv and Oren Tsur. 2021. From individuals | | 665 |
| | to communities : community-aware language | | 666 |
| | modeling for the detection of hate speech. | | 667 |
| | Eshwar Chandrasekharan, Umashanthi | | 668 |
| | Pavalanathan, Anirudh Srinivasan, Adam | | 669 |
| | Glynn, Jacob Eisenstein, and Eric Gilbert. 2017. | | 670 |
| | You can’t stay here: The efficacy of reddit’s 2015 | | 671 |
| | ban examined through hate speech. <i>Proceedings</i> | | 672 |
| | <i>of the ACM on Human-Computer Interaction</i> , | | 673 |
| | 1(CSCW):31. | | 674 |
| | Mithun Das, Punyajoy Saha, Ritam Dutt, Pawan | | 675 |
| | Goyal, Animesh Mukherjee, and Binny Mathew. | | 676 |
| | 2021. You too brutus! trapping hateful users in | | 677 |
| | social media: Challenges, solutions & insights. | | 678 |
| | In <i>Proceedings of the 32nd ACM Conference on Hy-</i> | | 679 |
| | <i>perertext and Social Media</i> , pages 79–89. | | 680 |
| | Thomas Davidson, Dana Warmley, Michael Macy, | | 681 |
| | and Ingmar Weber. 2017. Automated hate | | 682 |
| | speech detection and the problem of offensive | | 683 |
| | language. In <i>Eleventh international aai conference</i> | | 684 |
| | <i>on web and social media.</i> | | 685 |
| | Mai ElSherief, Caleb Ziem, David Muchlinski, | | 686 |
| | Vaishnavi Anupindi, Jordyn Seybolt, Munmun | | 687 |
| | De Choudhury, and Diyi Yang. 2021. Latent ha- | | 688 |
| | tred: A benchmark for understanding implicit | | 689 |
| | hate speech. <i>arXiv preprint arXiv:2109.05322.</i> | | 690 |
| | Paula Fortuna, Monica Dominguez, Leo Wanner, | | 691 |
| | and Zeerak Talat. 2022. Directions for NLP prac- | | 692 |
| | tices applied to online hate speech detection. In | | 693 |
| | <i>Proceedings of the 2022 Conference on Empirical</i> | | 694 |
| | <i>Methods in Natural Language Processing</i> , pages | | 695 |
| | 11794–11805, Abu Dhabi, United Arab Emirates. | | 696 |
| | Association for Computational Linguistics. | | 697 |
| | Linda Gordon. 2017. <i>The second coming of the KKK:</i> | | 698 |
| | <i>The Ku Klux Klan of the 1920s and the American</i> | | 699 |
| | <i>political tradition.</i> Liveright Publishing. | | 700 |
| | Jarod Govers, Philip Feldman, Aaron Dant, and | | 701 |
| | Panos Patros. 2023. Down the rabbit hole: De- | | 702 |
| | tecting online extremism, radicalisation, and | | 703 |
| | politicised hate speech. <i>ACM Computing Sur-</i> | | 704 |
| | <i>veys</i> , 55(14s):1–35. | | 705 |
| | William L. Hamilton, Rex Ying, and Jure Leskovec. | | 706 |
| | 2018. Inductive representation learning on large | | 707 |
| | graphs. <i>Preprint</i> , arXiv:1706.02216. | | 708 |
| | Abraham Israeli and Oren Tsur. 2022. Free speech | | 709 |
| | or free hate speech? analyzing the prolifera- | | 710 |
| | tion of hate speech in parler. In <i>Proceedings of</i> | | 711 |
| | <i>the Sixth Workshop on Online Abuse and Harms</i> | | 712 |
| | (WOAH), pages 109–121. | | 713 |
| | Urja Khurana, Ivar Vermeulen, Eric Nalisnick, Mar- | | 714 |
| | loes Van Noorloos, and Antske Fokkens. 2022. | | 715 |
| | Hate speech criteria: A modular approach to | | 716 |

| | | | |
|-----|--|---|---|
| 717 | task-specific hate speech definitions . In <i>Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)</i> , pages 176–191, Seattle, Washington (Hybrid). Association for Computational Linguistics. | Manoel Ribeiro, Pedro Calais, Yuri Santos, Virgílio Almeida, and Wagner Meira Jr. 2018. Characterizing and detecting hateful users on twitter. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 12. | 768 769 770 771 772 |
| 722 | Jaehoon Kim, Seungwan Jin, Sohyun Park, Someen Park, and Kyungsik Han. 2024. Label-aware hard negative sampling strategies with momentum contrastive learning for implicit hate speech detection . <i>Preprint</i> , arXiv:2406.07886. | Manoel Horta Ribeiro, Pedro H Calais, Yuri A Santos, Virgílio AF Almeida, and Wagner Meira Jr. 2017. "like sheep among wolves": Characterizing hateful users on twitter. <i>arXiv preprint arXiv:1801.00317</i> . | 773 774 775 776 777 |
| 727 | Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. <i>arXiv preprint arXiv:1609.02907</i> . | Paul Röttger, Haitham Seelawi, Debora Nozza, Zeerak Talat, and Bertie Vidgen. 2022. Multilingual hatecheck: Functional tests for multilingual hate speech detection models. <i>arXiv preprint arXiv:2206.09917</i> . | 778 779 780 781 782 |
| 730 | Lee Knuttila. 2011. User unknown: 4chan, anonymity and contingency. <i>First Monday</i> , 16(10). | Paul Röttger, Bertie Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet Pierrehumbert. 2021. HateCheck: Functional tests for hate speech detection models . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 41–58, Online. Association for Computational Linguistics. | 783 784 785 786 787 788 789 790 791 |
| 733 | Zachary Laub. 2019. Hate speech on social media: Global comparisons . | Haji Mohammad Saleem, Kelly P Dillon, Susan Benesch, and Derek Ruths. 2017. A web of hate: Tackling hateful speech in online social spaces. <i>arXiv preprint arXiv:1709.10159</i> . | 792 793 794 795 |
| 735 | Shu Li, Nayyar A Zaidi, Qingyun Liu, and Gang Li. 2021. Neighbours and kinsmen: hateful users detection with graph neural network. In <i>Advances in Knowledge Discovery and Data Mining: 25th Pacific-Asia Conference, PAKDD 2021, Virtual Event, May 11–14, 2021, Proceedings, Part I</i> , pages 434–446. Springer. | Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019a. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>ArXiv</i> , abs/1910.01108. | 796 797 798 799 |
| 742 | Simon Malevich and Tom Robertso. 2019. Violence begetting violence: An examination of extremist content on deep web social networks . | Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019b. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> . | 800 801 802 803 |
| 745 | Binny Mathew, Ritam Dutt, Pawan Goyal, and Animesh Mukherjee. 2019. Spread of hate speech in online social media. In <i>Proceedings of the 10th ACM conference on web science</i> , pages 173–182. | Zeerak Talat, James Thorne, and Joachim Bingel. 2018. Bridging the Gaps: Multi Task Learning for Domain Transfer of Hate Speech Detection , pages 29–55. Springer International Publishing, Cham. | 804 805 806 807 |
| 749 | Reid McIlroy-Young and Ashton Anderson. 2019. From “welcome new gabbers” to the pittsburgh synagogue shooting: The evolution of gab. In <i>Proceedings of the international aaii conference on web and social media</i> , volume 13, pages 651–654. | Kiran K. Thekumparampil, Chong Wang, Sewoong Oh, and Li-Jia Li. 2018. Attention-based graph neural network for semi-supervised learning . <i>Preprint</i> , arXiv:1803.03735. | 808 809 810 811 |
| 754 | David W McMillan and David M Chavis. 1986. Sense of community: A definition and theory. <i>Journal of community psychology</i> , 14(1):6–23. | Daniel Robert Thomas and Laila A Wahedi. 2023. Disrupting hate: The effect of deplatforming hate organizations on their online audience. <i>Proceedings of the National Academy of Sciences</i> , 120(24):e2214080120. | 812 813 814 815 816 |
| 757 | Zhenxiong Miao, Xingshu Chen, Haizhou Wang, Rui Tang, Zhou Yang, and Wenyi Tang. 2022. Detecting offensive language on social networks: An end-to-end detection method based on graph attention networks . <i>arXiv preprint</i> . | Elise Thomas. 2019. Aspi explains: 8chan . | 817 |
| 762 | Luke Munn. 2019. Alt-right pipeline: Individual journeys to extremism online . | Abhishek Velankar, Hrushikesh Patil, and Raviraj Joshi. 2022. A review of challenges in machine learning based automated hate speech detection . <i>Preprint</i> , arXiv:2209.05294. | 818 819 820 821 |

- 822 Zeerak Waseem and Dirk Hovy. 2016. Hateful sym-
823 bols or hateful people? predictive features for
824 hate speech detection on twitter. In *Proceedings*
825 *of the NAACL student research workshop*, pages
826 88–93.
- 827 Barry Wellman and Milena Gulia. 1999. Virtual
828 communities as communities. *Communities in*
829 *cyberspace*, pages 167–194.
- 830 Savvas Zannettou, Barry Bradlyn, Emiliano
831 De Cristofaro, Haewoon Kwak, Michael Siriv-
832 ianos, Gianluca Stringini, and Jeremy Blackburn.
833 2018. What is gab: A bastion of free speech or
834 an alt-right echo chamber. In *Companion Pro-*
835 *ceedings of the The Web Conference 2018*, pages
836 1007–1014. International World Wide Web Con-
837 ferences Steering Committee.
- 838 Caleb Ziems, William Held, Omar Shaikh, Jiaao
839 Chen, Zhehao Zhang, and Diyi Yang. 2024. [Can](#)
840 [Large Language Models Transform Computa-](#)
841 [tional Social Science?](#) *Computational Linguistics*,
842 50(1):237–291.

843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886

A Prompts and Generated Responses

Which of the following utterances may be perceived as hate speech?

1. Good morning boys and girls!
2. Muslim terrorists did this and Obama brought them here to run our government. Gotta get rid of Obammy, Omar, and all these other fucking Muslims.
3. Is there a more debased creature in our modern pop-cultural cesspool than (((Amy Schumer)))? Probably not.
4. We must secure the existence of our people and a future for white children!.
5. If white people are going to continue existing - we must support our race.
6. When the youngest white generation becomes adult they will be 5% of the world population, down from 1/3 in 1900.
7. Did you know that most dog breeds were bred by whites? We're awesome and they're awesome!
8. White lives matter!
9. Blue lives matter!
10. Black lives matter!

GPT-4 The utterances that may be perceived as hate speech are those that explicitly or implicitly promote hate, discrimination, or hostility towards particular groups based on their ethnicity, religion, or other inherent characteristics. Based on the examples you've provided, here's an analysis of which might be considered hate speech:

2. **"Muslim terrorists did this and Obama brought them here to run our government. Gotta get rid of Obammy, Omar, and all these other fucking Muslims"** - This statement is clearly hate speech. It uses derogatory language targeting Muslims and makes unfounded, harmful generalizations about an entire religious group.

3. **"Is there a more debased creature in our modern pop-cultural cesspool than (((Amy Schumer)))? Probably not."** - This statement might be perceived as hate speech, particularly because the use of triple parentheses is often a coded way to identify someone as Jewish in an anti-Semitic context. 887
888
889
890
891
892
893
894
4. **"We must secure the existence of our people and a future for white children!"** - This phrase is widely recognized as a white supremacist slogan, known as the "Fourteen Words." It is used to promote the idea of white racial superiority and is considered hate speech. 895
896
897
898
899
900
901
5. **"If white people are going to continue existing - we must support our race"** - This could be interpreted as promoting racial supremacy or separatism, which might be viewed as hate speech depending on the context in which it is used. 902
903
904
905
906
907

The other statements, while potentially controversial or promoting particular views, do not explicitly convey hate against a group based on inherent characteristics in the same way that the identified statements do. 908
909
910
911
912
913
914

Gemini 1.5 Pro Analysis of Potentially Hateful Utterances: Here's a breakdown of the provided utterances and their potential to be perceived as hate speech: 915
916
917
918

1. **Good morning boys and girls!** 919
 - This is a simple greeting and does not contain any hateful content. 920
921
2. **Muslim terrorists did this and Obama brought them here to run our government. Gotta get rid of Obammy, Omar and all these other fucking Muslims.** 922
 - This statement employs harmful stereotypes against Muslims, associating them with terrorism and expressing a desire to remove them from the government. This is a clear example of hate speech. 923
924
925
926
927
928
929
930
931
3. **Is there a more debased creature in our modern pop-cultural cesspool than (((Amy Schumer)))? Probably not.** 932
933
934

| | | | |
|-----|---|--|------|
| 935 | • The use of triple parentheses around a name is a known anti-semitic symbol online. This statement is likely intended to target Amy Schumer for her Jewish heritage and is therefore hate speech. | dismiss or oppose calls for police accountability and racial justice. | 981 |
| 936 | | | 982 |
| 937 | | • <i>Black lives matter</i> aims to highlight the systemic racism and violence faced by Black communities. It is a call for justice and equality, not a statement of superiority or hatred towards other groups. | 983 |
| 938 | | | 984 |
| 939 | | | 985 |
| 940 | | | 986 |
| 941 | 4. We must secure the existence of our people and a future for white children! | | 987 |
| 942 | | | 988 |
| 943 | 5. If white people are going to continue existing - we must support our race | <i>It's crucial to consider the context and intent behind these statements to determine if they constitute hate speech.</i> | 989 |
| 944 | | | 990 |
| 945 | 6. When the youngest white generation becomes adult they will be 5% of the world population, down from 1/3 in 1900 | | 991 |
| 946 | | B Networks Statistics | 992 |
| 947 | | Tables 6 and 7 provide more information on each of the datasets and justifies the focus on the largest connected component in each. | 993 |
| 948 | • These three statements, while not explicitly hateful, promote a white supremacist ideology. They focus on the preservation of white people and imply a fear of their decline, potentially inciting hostility towards other racial groups. | | 994 |
| 949 | | | 995 |
| 950 | | C Annotation Guidelines | 996 |
| 951 | | Each annotator was assigned 50 Parler accounts, ensuring that each account is being annotated by three students. For each account, annotators were given the user name, self-description of the user's account, and 30 posts published by the user. The 30 posts are sampled out of the following: (i) 15 posts with the highest hate speech prediction values ¹⁰ ; (ii) Its 10 most recent posts; (ii) Five random posts published by the account. The annotators showed a reasonable agreement level with an average Cohen's Kappa of 0.36. The full guidelines for the Parler dataset annotation process can be found here https://mega.nz/file/D0pjBDDR#9iMK0uC72w41BjSj_eDCmpbaJ6YGU8sZ8g0dvDv0pcU | 997 |
| 952 | | | 998 |
| 953 | | | 999 |
| 954 | | | 1000 |
| 955 | 7. Did you know that most dog breeds were bred by whites? We're awesome and they're awesome! | | 1001 |
| 956 | | | 1002 |
| 957 | | | 1003 |
| 958 | • This statement, while expressing pride in an accomplishment associated with white people, does not directly attack or demean any other group. It is unlikely to be considered hate speech. | | 1004 |
| 959 | | | 1005 |
| 960 | | | 1006 |
| 961 | | | 1007 |
| 962 | | | 1008 |
| 963 | | | 1009 |
| 964 | 8. White lives matter! | | 1010 |
| 965 | 9. Blue lives matter! | | 1011 |
| 966 | 10. Black lives matter! | | 1012 |
| 967 | • These statements, on the surface, simply affirm the value of different groups of people. However, depending on the context and intent, they can be interpreted in different ways. | D Responsible NLP Research Checklist | 1013 |
| 968 | | | 1014 |
| 969 | • <i>White lives matter</i> is often used as a response to the "Black Lives Matter" movement, implying a denial of the systemic racism faced by Black people. In such cases, it can be considered hate speech. | A1 Did you describe the limitations of your work? Yes | 1015 |
| 970 | | | 1016 |
| 971 | | A2 Did you discuss any potential risks of your work? Yes | 1017 |
| 972 | • <i>Blue lives matter</i> expresses support for law enforcement. While not inherently hateful, it can be used to | | 1018 |
| 973 | | A3 Do the abstract and introduction summarize the paper's main claims? Yes | 1019 |
| 974 | | | 1020 |
| 975 | | B Did you use or create scientific artifacts? Yes | 1021 |
| 976 | | | 1022 |
| 977 | | | |
| 978 | | | |
| 979 | | | |
| 980 | | | |

¹⁰The prediction is according to the BERT model introduced by Israeli and Tsur (2022)

| Dataset | Threshold | Precision | Recall | F1 | ROC AUC |
|---------|-----------|-----------|--------|-------|---------|
| Echo | 1 | 0.264 | 1.000 | 0.417 | 0.506 |
| | 3 | 0.266 | 1.000 | 0.420 | 0.511 |
| | 10 | 0.277 | 1.000 | 0.434 | 0.538 |
| | 50 | 0.358 | 0.950 | 0.520 | 0.673 |
| | 100 | 0.433 | 0.820 | 0.567 | 0.721 |
| Gab | 1 | 0.252 | 1.000 | 0.402 | 0.506 |
| | 3 | 0.257 | 1.000 | 0.409 | 0.519 |
| | 10 | 0.281 | 0.984 | 0.437 | 0.573 |
| | 50 | 0.380 | 0.559 | 0.452 | 0.628 |
| | 100 | 0.437 | 0.351 | 0.389 | 0.600 |
| Parler | 1 | 0.308 | 0.883 | 0.457 | 0.608 |
| | 3 | 0.341 | 0.824 | 0.482 | 0.644 |
| | 10 | 0.348 | 0.689 | 0.462 | 0.627 |
| | 50 | 0.457 | 0.387 | 0.420 | 0.616 |
| | 100 | 0.607 | 0.293 | 0.395 | 0.615 |

Table 5: Performance metrics for Echo, Gab, and Parler datasets

| Dataset | #Posts | #Users | #Edges | #Connected Components | #Singletons |
|--------------------|--------|--------|--------|-----------------------|-------------|
| Echo (Twitter) | 18M | 7.07K | 21.4K | 9075 | 2919 |
| Gab | 19.42M | 61.36K | 2.63M | 40K | 10.13K |
| Parler (this work) | 115M | 3.08M | 11.14M | 5.45M | 2.43M |

Table 6: Datasets statistics for the full network, excluding users without any posts.

| | | | | | |
|------|---|---|---|---|------|
| 1023 | B1 Did you cite the creators of artifacts you used? Yes | B5 Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.? Yes | 1046 | | |
| 1024 | | | 1047 | | |
| 1025 | B2 Did you discuss the license or terms for use and / or distribution of any artifacts? Yes, we discuss the limitations and ethical considerations of our work | B6 Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Yes | 1048 | | |
| 1026 | | | 1049 | | |
| 1027 | | | 1050 | | |
| 1028 | | | 1051 | | |
| 1029 | B3 Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)? Yes | C Did you run computational experiments? Yes | 1052 | | |
| 1030 | | | 1053 | | |
| 1031 | | | C1 Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? Yes | 1054 | |
| 1032 | | | | 1055 | |
| 1033 | | | | 1056 | |
| 1034 | | | B4 Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it? No, since our work involves harmful or offensive content. | C2 Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? Yes | 1057 |
| 1035 | | | | | 1058 |
| 1036 | 1059 | | | | |
| 1037 | 1060 | | | | |
| 1038 | | C3 Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether | 1061 | | |
| 1039 | | | 1062 | | |
| 1040 | | | 1063 | | |
| 1041 | | | 1064 | | |
| 1042 | | | 1065 | | |
| 1043 | | | 1066 | | |
| 1044 | | | | | |
| 1045 | | | | | |

| Dataset | #Posts | #Users | #Edges | Raw Data | | Annotated Data | |
|--------------------|--------|--------|--------|------------------------|---------------|----------------|--------|
| | | | | Clustering Coefficient | Optimal Gamma | #Users | % Hate |
| Echo (Twitter) | 9.8M | 3.7K | 20.7K | 0.19 | 2.8 | 532 | 26.1% |
| Gab | 19.28M | 51.2K | 2.47M | 0.402 | 4.06 | 982 | 24.5% |
| Parler (this work) | 60.7M | 643K | 11.4M | 0.224 | 2.14 | 881 | 25.2% |

Table 7: Datasets statistics for the largest (weakly) connected component. We treated the graphs as undirected in order to calculate the Clustering Coefficient and the Optimal Gamma (assuming a Power-law distribution)

1067 you are reporting the max, mean, etc. or
1068 just a single run? **Yes**

1069 C4 If you used existing packages (e.g., for
1070 preprocessing, for normalization, or for
1071 evaluation), did you report the imple-
1072 mentation, model, and parameter settings
1073 used (e.g., NLTK, Spacy, ROUGE, etc.)?
1074 **Yes**

1075 D Did you use human annotators (e.g.,
1076 crowdworkers) or research with human
1077 participants? **Yes**

1078 D1 Did you report the full text of instruc-
1079 tions given to participants, including e.g.,
1080 screenshots, disclaimers of any risks to
1081 participants or annotators, etc.? **Yes**

1082 D2 Did you report the full text of instruc-
1083 tions given to participants, including e.g.,
1084 screenshots, disclaimers of any risks to
1085 participants or annotators, etc.? **Yes**

1086 D3 Did you discuss whether and how con-
1087 sent was obtained from people whose
1088 data you’re using/curating **Yes**

1089 D4 Was the data collection protocol approved
1090 (or determined exempt) by an ethics re-
1091 view board? **N/A**

1092 D5 Did you report the basic demographic
1093 and geographic characteristics of the an-
1094 notator population that is the source of
1095 the data? **Yes**

1096 E Did you use AI assistants (e.g., ChatGPT,
1097 Copilot) in your research, coding, or writ-
1098 ing? **Yes [E1]** Did you include informa-
1099 tion about your use of AI assistants? **We**
1100 **used LLMs as classifiers, see Table 1 and**
1101 **references through the paper.**