

# Hate Speech Detection with Generalizable Target-aware Fairness

Tong Chen

The University of Queensland  
Brisbane, Australia  
tong.chen@uq.edu.au

Danny Wang

The University of Queensland  
Brisbane, Australia  
danny.wang@uq.edu.au

Xurong Liang

The University of Queensland  
Brisbane, Australia  
xurong.liang@uq.edu.au

Marten Risius

The University of Queensland  
Brisbane, Australia  
m.risius@business.uq.edu.au

Gianluca Demartini

The University of Queensland  
Brisbane, Australia  
g.demartini@uq.edu.au

Hongzhi Yin\*

The University of Queensland  
Brisbane, Australia  
h.yin1@uq.edu.au

## Abstract

To counter the side effect brought by the proliferation of social media platforms, hate speech detection (HSD) plays a vital role in halting the dissemination of toxic online posts at an early stage. However, given the ubiquitous topical communities on social media, a trained HSD classifier can easily become biased towards specific targeted groups (e.g., *female* and *black* people), where a high rate of either false positive or false negative results can significantly impair public trust in the fairness of content moderation mechanisms, and eventually harm the diversity of online society. Although existing fairness-aware HSD methods can smooth out some discrepancies across targeted groups, they are mostly specific to a narrow selection of targets that are assumed to be known and fixed. This inevitably prevents those methods from generalizing to real-world use cases where new targeted groups constantly emerge (e.g., new forums created on Reddit) over time. To tackle the defects of existing HSD practices, we propose Generalizable target-aware Fairness (GetFair), a new method for fairly classifying each post that contains diverse and even unseen targets during inference. To remove the HSD classifier’s spurious dependence on target-related features, GetFair trains a series of filter functions in an adversarial pipeline, so as to deceive the discriminator that recovers the targeted group from filtered post embeddings. To maintain scalability and generalizability, we innovatively parameterize all filter functions via a hypernetwork. Taking a target’s pretrained word embedding as input, the hypernetwork generates the weights used by each target-specific filter on-the-fly without storing dedicated filter parameters. In addition, a novel semantic gap alignment scheme is imposed on the generation process, such that the produced filter function for an unseen target is rectified by its semantic affinity with existing targets used for training. Finally, experiments<sup>1</sup> are conducted on

two benchmark HSD datasets, showing advantageous performance of GetFair on out-of-sample targets among baselines.

## CCS Concepts

• **Information systems** → **Web mining; Data mining; Security and privacy** → *Social aspects of security and privacy.*

## Keywords

Hate Speech Detection; Target-aware Fairness; Debaised Content Moderation; Data Science for Social Good

## ACM Reference Format:

Tong Chen, Danny Wang, Xurong Liang, Marten Risius, Gianluca Demartini, and Hongzhi Yin. 2024. Hate Speech Detection with Generalizable Target-aware Fairness. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD ’24)*, August 25–29, 2024, Barcelona, Spain. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3637528.3671821>

## 1 Introduction

Many benefits of social media’s liberation of communication come at the expense of proliferating hate speech. To prevent the negative socio-economical impact from hateful content, accurate algorithms for hate speech detection (HSD) have been heavily investigated by both industry practitioners [12] and research communities [14].

Meanwhile, on the flip side of the coin, it has been reported [15] that various HSD algorithms have exposed vulnerability to different types of biases, including identity, annotation, and political biases. In the context of HSD, these biases are predominately associated with the sources (i.e., authors) of online posts, and a variety of corresponding solutions have also been made available. However, there also exists bias towards the *targets*, or more precisely the posts’ *targeted groups* [15], which are usually an identity group (e.g., African Americans), or a particular protected user attribute (e.g., one’s religion). In the rest of this paper, when there is no ambiguity, we will use target to refer to targeted groups in a post. The targets of a post can be identified based on the main topics discussed, the context of conversation, or the channels hosting this post (e.g., the *incel* forum on Reddit). Due to the inherent disparity in label distributions and language styles among different targets, models trained on such skewed data can reflect highly unstable HSD performance across targets. As a result, it is commonly seen that HSD classifiers exhibit abnormally high false positive or negative rates on some targeted groups. A high false positive rate on a specific target means that, the HSD classifier is prone to misclassifying a neutral post as hateful

\*Corresponding author.

<sup>1</sup>The implementation of GetFair is released at <https://github.com/xurong-liang/GetFair>.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*KDD ’24, August 25–29, 2024, Barcelona, Spain*

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0490-1/24/08

<https://doi.org/10.1145/3637528.3671821>

as long as this target is mentioned, introducing a greater chance of seeing incorrectly blocked posts in social media applications. In contrast, the high rate of false negative predictions are a result of failing to recognize and moderate truly hateful content, which in turn makes the specific target group more vulnerable to being attacked online.

As such, achieving target-aware fairness in HSD becomes an increasingly important requirement for today’s content moderation practices, and an emerging area to be researched. The uneven target-wise HSD performance implies that a spurious correlation between specific targets and the labels [30, 44] (e.g., *hateful* and *neutral* in a binary setting) is established by the classifier. The spurious correlation then misleads the HSD classifier to provide incorrect labels with a groundlessly high confidence. To prevent a classifier from such biases, the ultimate goal is to lower its sensitivity to spurious features like target-related texts (i.e., *who* are being discussed), and instead base its judgement on more generalizable linguistic features (i.e., *how* they are being discussed). In this way, balanced HSD performance across targets can be attained.

To pursue this goal, relevant literature has seen various data-centric methods that append new data sources or alter the way of engaging data samples to debias HSD models. Examples include sample reweighting [48, 69, 72] that weakens the weight of samples with spurious features and lays more emphasis on samples that are less prone to biases during training, as well as marking bias-sensitive words [12, 21, 44] to perform label corrections or regulate model-level fairness. Nevertheless, these solutions mostly come with a strong empirical nature, such as the need for identifying confounding lexical patterns beforehand [48] and the labor-intensive process of manual annotation [44]. This has stimulated the emergence of less empirical, model-centric fair HSD solutions that are often designed under the principle of “fairness through unawareness” [36]. In a nutshell, filtering modules are designed as an add-on to the HSD model, where information related to protected attributes are implicitly removed from learned representations of an online post. Then, methods like regularization [19], multi-task learning [18], and adversarial training [26] are adopted to jointly optimize the information filters along with the HSD task.

Despite the stronger performance of filter-based solutions and their plug-and-play compatibility with most HSD classifiers, their practicality is inevitably challenged by the diversity of targets in real-world HSD applications. Generally, the majority of the target debiasing methods in HSD bear an assumption that the number and identity of targets are fixed and consistent (e.g., only *race* and *gender*) across both training and inference phases. However, due to the difficulty in obtaining high-quality labeled datasets for HSD [15], the targets of acquired posts in the training data are just a tip of the iceberg compared with the ones in the online environment, rendering this assumption ill-posed and hurting the generalizability of existing debiasing approaches. For most solutions, their filters are only designed and trained for targets that are known in the training stage, and are thus hardly transferrable to an unseen target group. Furthermore, a single filter function, commonly parameterized as a neural network such as multi-layer perceptron (MLP), is subject to limited capacity for removing multifaceted target-specific information from textual embeddings, bringing the need for target-specific content filters. Considering the amount of possible targets being

discussed online, it is infeasible to learn a dedicated content filter per target from both scalability and data availability perspectives. Given the high velocity of evolving targets, and naturally emerging new targets of posts on social media platforms (e.g., new Reddit forums and Facebook hashtags created daily), the capability of generalizing to new targets without being constantly retrained is deemed crucial.

To this end, in this paper we propose a novel approach, namely HSD with Generalizable target-aware Fairness (GetFair). In GetFair, to avoid the linear parameter cost associated with the number of target-specific filter functions, instead of independently parameterizing each filter, we put forward a hypernetwork [63] that adaptively generates filter parameters – MLP weights and biases in our case – for each target. When debiasing post embeddings with an identified target described by one or several keywords, the input to the hypernetwork is defined as the target’s indicator embedding, which can be easily composed with off-the-shelf pretrained word embeddings<sup>2</sup>. On the one hand, as the target embedding serves as the indicator to control the parameter generation process, the uniqueness of generated target-specific filters is guaranteed. On the other hand, by ensembling target-specific filters, GetFair is able to efficiently handle cases where a combination of different targets are identified in the same post, which is a trickier yet understudied case in real-world scenarios [27]. Moreover, the readily available target embeddings makes it possible for the hypernetwork to generate filter parameters with a valid indicator for an arbitrary, unseen target. An adversarial training paradigm is in place to optimize the filters toward learning debiased post embeddings. Concretely, a target classifier is deployed as the discriminator and tries to infer the original targets from the filtered post embeddings, where a capable filter function is ultimately learned to maximize the classification error of the discriminator. To further refine the dynamic filter generation for each target that arrives, a semantic gap alignment scheme is proposed, such that the semantic distance among all targets’ input embeddings is resembled in the parameter space of their corresponding filters.

To sum up, the contributions of this paper are three-fold:

- We focus on an emerging problem in HSD concerning prediction fairness across different targeted groups of online posts by isolating the HSD classifier’s judgement from a post’s target-related spurious features. We point out the generalization deficiency of existing target-aware HSD debiasing methods when facing newly identified targets that are unknown during training.
- A remedy, namely GetFair is proposed to achieve generalizable target-aware fairness as a plug-and-play method for most HSD classifiers. Instead of training a filter function for each target, GetFair optimizes a hypernetwork that adaptively parameterizes target-specific filters based on the indicators it receives. The generated filters are adversarially trained and further regularized via an innovative semantic gap alignment constraint, uplifting their debiasing capability without the pressure on scalability.
- We conduct extensive experiments on two public benchmark datasets, where the results have demonstrated the superiority

<sup>2</sup>GetFair adopts GloVe [42] word embeddings.

of GetFair, evidenced by its balanced effectiveness-fairness trade-off in HSD tasks on out-of-sample targets compared with state-of-the-art HSD debiasing baselines.

## 2 Preliminaries

In this section, we mathematically define the concept of target-aware fairness in HSD, metrics for quantifying such fairness, and our research objective w.r.t. generalizable target-aware HSD.

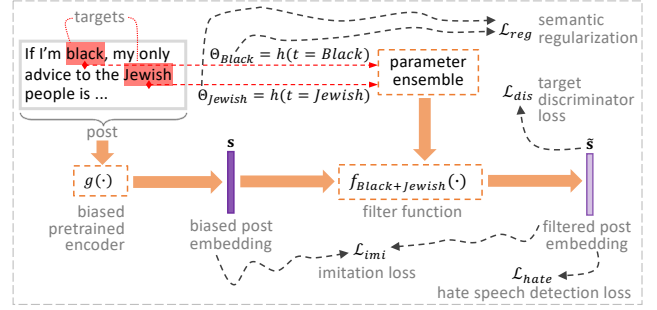
**Hate Speech Detection Tasks.** Hate speech detection is commonly defined as a classification task. In our paper, the default task setting is to take a social media post  $s$  consisting of a sequence of tokens as the input, and output a predicted scalar  $\hat{y} \in (0, 1)$  to represent the likelihood of  $s$  being a hate speech. In the most recent literature [35, 59], pretrained transformer models like BERT [62] and GPT [43] are commonly the default encoder  $g(\cdot)$  for generating the post-level embeddings  $\mathbf{s} = g(s)$  given their superiority in representation learning, which is also the initial feature we feed into the HSD classifier. Note that our paper is scoped around this text-based, binary classification setting, which is the most widely used one [15] in HSD. Meanwhile, with the availability of more nuanced annotated data, our findings can be easily generalized to some emerging settings like multi-class classification with different levels of hatefulness [68] or HSD with additional multimodal data (e.g., images [24]) given their text-focused origin.

**Definition 1: Targets in A Post.** In our setting, each social media post  $s_i$  has one or more targeted user groups being discussed. In  $s_i$ , those mentioned user groups are termed *targets*, denoted by  $t \in \mathcal{T}_i$ . If deployed in a real online environment, one can consider using rule- or lexical-based approaches [13, 55], or trainable detectors [52] to identify the target set  $\mathcal{T}_i$  from each post  $s_i$ . Because our emphasis is to uplift target-aware HSD fairness, we directly employ datasets (see Section 4.1 for details) that come with target labels, which are identified by human annotators from the textual clues within each post [35, 41, 54]. The target identification task, as a more widely studied task, is not investigated in the paper.

**Definition 2: Target-aware Fairness Metrics.** We assume a finite set of targeted groups  $\mathcal{T}$  used for evaluation. For posts containing each target  $t \in \mathcal{T}$ , the overall group-level HSD classification performance can be measured by a specific accuracy metric. Then, fairness across targeted groups is reflected by the performance discrepancy among all  $|\mathcal{T}|$  targets. Intuitively, the smaller the target-wise discrepancy, the fairer the HSD classifier. In our work, to quantify the level of algorithmic fairness, we adopt the well-established notion of *False Positive/Negative Equality Difference* [12] suggested by Google Jigsaw, abbreviated as FPED/FNED. The only slight difference is that we normalize their values with the number of targets involved to make them comparable across different datasets, and the normalized versions are termed nFPED and nFNED:

$$\begin{aligned} \text{nFPED} &= \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} |\text{FPR} - \text{FPR}(t)|, \\ \text{nFNED} &= \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} |\text{FNR} - \text{FNR}(t)|, \end{aligned} \quad (1)$$

where  $\text{FPR}(t)$  and  $\text{FNR}(t)$  are respectively the false positive and false negative rates on posts associated with target  $t$ , while FPR and FNR respectively denote the overall false positive and false



**Figure 1: An overarching view of GetFair. Detailed designs of the four objective functions can be found in Sections 3.2 ( $\mathcal{L}_{reg}$ ), 3.3 ( $\mathcal{L}_{dis}$ ), and 3.4 ( $\mathcal{L}_{hate}$  and  $\mathcal{L}_{imi}$ ), respectively.**

negative rates on all posts. Essentially, both nFPED and nFNED measure how far on average the HSD performance on each target  $t$  deviates from the global average, and ideally will converge to zero if all targets are treated evenly by the HSD classifier. Since it is equally important to account for both nFPED and nFNED for fairness, we further add the harmonic mean of both metrics, termed *harmonic fairness* (HF):

$$\text{HF} = 2 \times \frac{\text{nFPED} \times \text{nFNED}}{\text{nFPED} + \text{nFNED}}, \quad (2)$$

which jointly factors in two metrics for fairness evaluation. With the fairness metrics defined, we formulate our HSD task with generalizable target-aware fairness below.

**Definition 3: Hate Speech Detection with Generalizable Target-aware Fairness.** Given a set of posts  $\mathcal{S}_{train} = \{(s_i, \mathcal{T}_i, y_i)\}_{i=1}^{|\mathcal{S}_{train}|}$  where the  $i$ -th post  $s_i$  is associated with the binary label  $y_i$  and a subset of known targets  $\mathcal{T}_i \in \mathcal{T}_{train}$ ,  $\mathcal{T}_i \neq \emptyset$ , our objective is to train a debiased HSD classifier  $\phi_{\text{HSD}}(\cdot)$ . The trained classifier  $\phi_{\text{HSD}}(\cdot)$  is evaluated on the test set  $\mathcal{S}_{test} = \{(s_j, \mathcal{T}_j, y_j)\}_{j=1}^{|\mathcal{S}_{test}|}$ , where  $\mathcal{T}_j \in \mathcal{T}_{test}$ ,  $\mathcal{T}_j \neq \emptyset$ . More importantly, all test-time posts will contain at least one target that is *completely unseen* during training, i.e.,  $\mathcal{T}_{test} \setminus \mathcal{T}_{train} \neq \emptyset$ . Ideally,  $\phi_{\text{HSD}}(\cdot)$  is expected to: (1) provide maximal classification effectiveness, reflected by larger scores in accuracy, F1, and AUC; and (2) yield minimal target-wise bias, reflected by lower scores in nFPED, nFNED, and HF.

## 3 GetFair: Model Design

We provide a graphical view of GetFair’s workflow in Figure 1. With our research task defined, we unfold the design of each core component of it.

### 3.1 Target-specific Filter Generation with Adaptive Hypernetwork

In GetFair, we utilize a filter function  $f(\cdot)$  which takes the generated  $d$ -dimensional content embedding  $\mathbf{s} \in \mathbb{R}^d$  from an arbitrary encoder  $g(\cdot)$  as its input, and emits a debiased representation  $\tilde{\mathbf{s}} = f(\mathbf{s}) \in \mathbb{R}^d$ . The filter  $f(\cdot)$  will be trained to remove target-specific information from the content embedding  $\mathbf{s}$ , such that subsequent hate speech predictions made with  $\tilde{\mathbf{s}}$  are less reliant on it and will exhibit less performance bias on different targets while maintaining accuracy. Without loss of generality,  $f(\cdot)$  can be formulated as a multi-layer perceptron (MLP) with weights  $\mathbf{W} \in \mathbb{R}^{d \times d}$  and biases  $\mathbf{b} \in \mathbb{R}^d$  at

all layers. Considering the vast pool of possible targets in online platforms, it is unrealistic to rely on a single MLP’s capacity for removing all targets’ relevant information from  $\mathbf{s}$ . Thus, a more performant approach is to instantiate one filter function  $f_t(\cdot)$  per target  $t$ . However, this will inevitably lead to difficulties in scaling to numerous targets in reality and generalizing to unseen targets (e.g., an emerging buzzword) during inference.

To ensure the practicality of GetFair as a plug-in debiasing method for HSD, we propose an efficient alternative for building target-aware filter functions. Instead of letting every filter function have a dedicated set of trainable parameters, we leverage a hypernetwork [63]  $h(\cdot)$  that can adaptively parameterize each target-specific MLP filter. Concretely, for every target  $t$ , its dedicated filter parameters are generated via:

$$\Theta_t^{(l)} = [\mathbf{W}_t^{(l)}, \mathbf{b}_t^{(l)}] = h_l(\mathbf{t}), \quad (3)$$

where  $h_l(\cdot)$  is the hypernetwork responsible for generating parameters (i.e., weights and biases) for the  $l$ -th layer of filter  $f_t(\cdot)$ , and  $[\cdot, \cdot]$  denotes horizontal concatenation. For notation simplicity, we use  $\Theta_t^{(l)} \in \mathbb{R}^{d \times (d+1)}$  to denote the concatenated weight and bias at each layer. Following common practices in hypernetworks [51, 64],  $h_l(\cdot)$  firstly outputs a flat vector with  $d^2 + d$  dimensionality, which is further reshaped into the matrix form of  $\Theta_t^{(l)}$ .

**Target Indicators.** Notably,  $h_l(\cdot)$  is shared across all targets, and the generated parameters are solely conditioned on its input  $\mathbf{t} \in \mathbb{R}^{d'}$ , which is the target indicator that informs the hypernetwork of the exact target filter to generate. To facilitate effective information filtering and target-wise generalizability, the designed target indicators should meet two characteristics: (1) they should be distinguishable for different targets such that every target-specific filter maintains uniqueness; and (2) they are capable of representing an arbitrary number of unseen targets without any training. Hence, this eliminates some indicator methods commonly used in hypernetworks such as one-hot encoding and learnable embeddings [51, 63, 64]. In GetFair, we take advantage of the wide availability of pretrained word embeddings, namely the GloVe embeddings [42] for composing the target indicator  $\mathbf{t}$ . We have adopted its 300-dimensional version, thus  $d' = 300$ . Specifically, for target  $t = \{v_1, \dots, v_k, \dots, v_n\}$  represented as a sequence of  $n$  vocabularies, we take the mean of all  $n$  corresponding word embeddings  $\mathbf{v}_k \in \mathbb{R}^{d'}$  as its representation, i.e.,  $\mathbf{t} = \frac{1}{n} \sum_{k=1}^n \mathbf{v}_k$ .

**Low-rank Parameterization.** In the default design of the adaptive hypernetwork, the output space for each hypernetwork is  $d \times (d + 1)$ , which is substantially large (e.g., 16, 512 predictions to make for  $d = 128$ ) compared with the input dimensionality. Consequently, this creates a low-dimensional bottleneck where the hypernetwork layers are not sufficiently expressive for the downstream predictions [4, 17, 66], thus being prone to underfitting. Furthermore, although the target-specific filter parameters are no longer stored owing to the adaptive hypernetwork, they are still required for in-memory forward and backward passes during run time, thus significantly limiting the batch size allowed and negatively impacting the training time and memory efficiency. As such, we adopt a low-rank formulation for the filter parameters:

$$\Theta_t^{(l)} = \mathbf{U}_t^{(l)} \mathbf{W}_t^{(l)} \mathbf{V}_t^{(l)}, \quad (4)$$

where  $\mathbf{U}_t^{(l)} \in \mathbb{R}^{d \times K}$ ,  $\mathbf{W}_t^{(l)} \in \mathbb{R}^{K \times K}$ ,  $\mathbf{V}_t^{(l)} \in \mathbb{R}^{K \times (d+1)}$ , and the rank  $K \ll d$ . With this low-rank parameterization, assuming  $d = d'$ , the memory cost for an  $L$ -layer target-specific filter is dramatically reduced from  $\mathcal{O}(Ld^2 + Ld)$  to  $\mathcal{O}(LK^2 + 2LdK + LK)$ , which now conveniently supports parallelized batch computing.

**Multi-target Filters.** It is worth mentioning that, in real applications, one post might be attributed to several different targets, calling upon the necessity for a mechanism that can simultaneously filter a set of targets  $\mathcal{T}_i$  associated with each post’s embedding  $\mathbf{s}_i$ . A straightforward way of dosing so is to remove one target-specific information at a time via  $\tilde{\mathbf{s}}_{i,t} = f_t(\mathbf{s}_i)$  for every  $t \in \mathcal{T}_i$ , then merge all filtered embeddings of the same post (e.g., via sum pooling [29]) into  $\tilde{\mathbf{s}}_i$ . Though each  $\tilde{\mathbf{s}}_{i,t}$  removes information related to target  $t \in \mathcal{T}_i$ , it is not necessarily the case for a different  $\tilde{\mathbf{s}}_{i,t'}$  where  $t' \in \mathcal{T}_i$  and  $t' \neq t$ , as the associated filter  $f_{t'}(\cdot)$  is only dedicated to removing information about  $t'$ . Consequently, the embedding fusion step may introduce unwanted target-specific information back into  $\tilde{\mathbf{s}}_i$ . To bypass this potential defect, we design a combinatorial approach to form a multi-target filter via parameter ensemble:

$$\tilde{\mathbf{s}}_i = f_{\mathcal{T}_i}(\mathbf{s}_i | \Theta_{\mathcal{T}_i}), \quad \Theta_{\mathcal{T}_i} = \{\Theta_t^{(l)}\}_{t \in \mathcal{T}_i}^L, \quad (5)$$

and the  $l$ -th layer of this multi-target filter  $f_{\mathcal{T}_i}(\cdot)$  is parameterized by  $\Theta_{\mathcal{T}_i}^{(l)} \in \Theta_{\mathcal{T}_i}$ , which is calculated as follows:

$$\Theta_{\mathcal{T}_i}^{(l)} = \frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} \Theta_t^{(l)}. \quad (6)$$

In a nutshell, given a post embedding  $\mathbf{s}_i$  and its targets  $\mathcal{T}_i$ , our combinatorial approach aggregates all target-specific filters’ parameters first, and then compute the filtered embedding in one go. As such, we can prevent the resultant post embedding  $\tilde{\mathbf{s}}_i$  from being contaminated by the late fusion of individually filtered embeddings.

### 3.2 Regularizing Filter Parameters with Semantic Gap Alignment

Compared with learning multiple filter functions, learning a unified adaptive hypernetwork brings substantially lighter parameterization, but also comes at a risk of overfitting some specific target distributions, e.g., the frequent ones, challenging the generalizability. Hence, the hypernetwork, if trained without regularization, can potentially fail to generalize when generating filter weights for completely unseen targets. To prevent the generated filter parameters from undesired homogeneity and thus impaired utility, we impose an additional regularization on the distribution of the generated filter parameters  $\Theta_t^{(l)}$  w.r.t. different targets. Specifically, for every pair of targets  $t, t'$  and their indicators  $\mathbf{t}, \mathbf{t}'$ , we would like to preserve their semantic distance within the generated filter parameters  $\Theta_t, \Theta_{t'}$  too, i.e.,  $d(\mathbf{t}, \mathbf{t}') \approx d(\Theta_t^{(l)}, \Theta_{t'}^{(l)})$  with a distance metric  $d(\cdot, \cdot)$ . To bypass the unmatched dimensionality between  $\mathbf{t}$  and  $\Theta_t$ , we define a semantic gap alignment scheme as the regularization loss  $\mathcal{L}_{reg}$  as follows:

$$\mathcal{L}_{reg} = \sum_{l=1}^L \sum_{t, t' \in \mathcal{T}_{train}} \left( \cos(\mathbf{t}, \mathbf{t}') - \cos(\bar{\Theta}_t^{(l)}, \bar{\Theta}_{t'}^{(l)}) \right)^2, \quad (7)$$

where  $\cos(\cdot, \cdot)$  is the cosine similarity, and  $\bar{\Theta} \in \mathbb{R}^{d(d+1)}$  is the flattened version of  $\Theta$  in vector format to facilitate cosine similarity

comparison between filter parameters. By imposing  $\mathcal{L}_{reg}$ , any two similar/dissimilar targets (as reflected by their indicators) will be assigned corresponding filter functions that reflect the same level of relatedness, assuring the semantic gaps among targets to be mirrored into the parameter space.

### 3.3 Target Discriminator

Ideally, each debiased post embedding  $\tilde{s}_i$  no longer carries target-specific information that leads to the hypothetically unfair classification results. To ensure that the target filter removes the target-related information, we propose to take advantage of adversarial training to optimize each filter. Specifically, we set up a target discriminator defined as the following:

$$\hat{\mathbf{p}}_i = \text{sigmoid}(\text{MLP}_{dis}(\tilde{s}_i)), \quad (8)$$

where  $\text{sigmoid}(\text{MLP}_{dis}(\cdot)) : \mathbb{R}^d \rightarrow (0, 1)^{|\mathcal{T}_{train}|}$  is an MLP with a sigmoid readout that maps the input  $\tilde{s}$  into a  $|\mathcal{T}_{train}|$ -dimensional output, with each entry representing the probability of having a target associated with the post  $\tilde{s}_i$ . Note that given the existence of multiple targets in a single post,  $\hat{\mathbf{p}}_i$  is essentially a collection of individual binary classification scores instead of a softmax distribution. The optimization of  $\text{MLP}_{dis}(\cdot)$  is facilitated by minimizing the point-wise log loss below:

$$\mathcal{L}_{dis} = - \sum_{\forall s_i \in \mathcal{S}_{train}} \left( \mathbf{p}_i^\top \log \hat{\mathbf{p}}_i + (1 - \mathbf{p}_i)^\top \log(1 - \hat{\mathbf{p}}_i) \right), \quad (9)$$

where  $\mathbf{p}_i \in \{0, 1\}^{|\mathcal{T}_{train}|}$  is the ground truth multi-hot label derived from every  $\mathcal{T}_i$ . We will defer the introduction of the entire adversarial training procedure to Section 3.5.

### 3.4 Hate Speech Classifier

As the debiasing objective is designed in parallel to the pure HSD task, we hereby define the classifier used for HSD. Simply put, the hate speech detector inherits the traditional binary classifier formulation in most existing HSD methods with its input updated from the biased post embedding  $s_i$  to the filtered one  $\tilde{s}_i$ :

$$\hat{y}_i = \text{sigmoid}(\text{MLP}_{hate}(\tilde{s}_i)), \quad (10)$$

where  $\text{sigmoid}(\text{MLP}_{hate}(\cdot)) : \mathbb{R}^d \rightarrow (0, 1)$  produces a probability scalar  $\hat{y}_i$  indicating the likelihood of  $\tilde{s}_i$  being a hateful post. Being a binary classification task, log loss is adopted for training:

$$\mathcal{L}_{hate} = \sum_{\forall s_i \in \mathcal{S}_{train}} \left( y_i \log \hat{y}_i + (1 - y_i) \log \hat{y}_i \right), \quad (11)$$

where  $y_i \in \{0, 1\}$  is the binary label of each training instance.

**Enhancing HSD Performance via Imitation Learning.** To prevent the filtered post embedding  $\tilde{s}$  from the byproduct of performance degradation, we further enhance the training of the HSD classifier with imitation learning. Specifically, we ask  $\text{MLP}_{hate}(\tilde{s}_i)$  to mimic its response when given the unfiltered version of the post embedding  $s_i$ . This is achieved by aligning the predicted probability distributions over the two binary classes:

$$\mathcal{L}_{imi} = D_{KL}([\hat{y}_i, 1 - \hat{y}_i] \parallel [\hat{y}'_i, 1 - \hat{y}'_i]), \quad (12)$$

where  $D_{KL}(\cdot \parallel \cdot)$  measures the Kullback-Leibler (KL) divergence between two distributions, and  $\hat{y}'_i = \text{sigmoid}(\text{MLP}_{hate}(s_i))$  is the prediction based on the unfiltered post embedding  $s_i$ . The imitation

---

#### Algorithm 1 Optimization Procedure of GetFair

---

```

1: Input:  $\mathcal{T}_{train}$ , parameters  $\Theta_{enc}$ ,  $\Theta_{hyper}$ ,  $\Theta_{dis}$ ,  $\Theta_{hate}$  respectively
   for the pretrained encoder, hypernetwork, target discriminator,
   and HSD classifier, epoch numbers  $N$  and  $N'$ , loss coefficients
    $\mu$ ,  $\gamma$ , and  $\lambda$ 
2: Output: Optimized  $\Theta_{enc}$ ,  $\Theta_{hyper}$ ,  $\Theta_{dis}$ ,  $\Theta_{hate}$ 
3: Randomly initialize  $\Theta_{hyper}$ ,  $\Theta_{dis}$ ,  $\Theta_{hate}$ ;
4: repeat
5:   for  $epoch_d = 1, \dots, K$  do
6:     Draw a mini-batch from  $\mathcal{S}_{train}$ ;
7:     Freeze  $\Theta_{hyper}$ ,  $\Theta_{hate}$  and  $\Theta_{enc}$ ;
8:      $\tilde{s} \leftarrow \text{Eq.}(5)$ ;
9:      $\mathcal{L}_{dis} \leftarrow \text{Eq.}(9)$ ;
10:    Update  $\Theta_{dis}$  w.r.t.  $\mathcal{L}_{dis}$ ;
11:   for  $epoch_f = 1, \dots, K'$  do
12:     Freeze  $\Theta_{dis}$ ;
13:     Draw a mini-batch from  $\mathcal{S}_{train}$ ;
14:      $\tilde{s} \leftarrow \text{Eq.}(5)$ ;
15:      $\mathcal{L}_{reg} \leftarrow \text{Eq.}(7)$ ,  $\mathcal{L}_{dis} \leftarrow \text{Eq.}(9)$ ,
        $\mathcal{L}_{hate} \leftarrow \text{Eq.}(10)$ ,  $\mathcal{L}_{imi} \leftarrow \text{Eq.}(12)$ ;
16:     Update  $\Theta_{hyper}$ ,  $\Theta_{hate}$  and  $\Theta_{enc}$  w.r.t.  $\mathcal{L}_{hate} + \mu \mathcal{L}_{reg}$ 
        $+ \gamma \mathcal{L}_{imi} - \lambda \mathcal{L}_{dis}$ ;
17: until convergence

```

---

loss  $\mathcal{L}_{imi}$  encourages the same hate speech classifier to emit similar decisions no matter whether the post embedding is filtered or not, thus further decorrelating the HSD classifier with the spurious target-related features in the texts.

### 3.5 Adversarial Optimization via Alternation

Alongside the fundamental goal of achieving accurate HSD, GetFair encompasses two additional components that have adversarial goals. The target discriminator is designed to identify associated targets for each given post embedding, while the filter function  $f(\cdot)$  essentially tries to deceive it such that relevant targets cannot be confidently inferred from the debiased embedding  $\tilde{s}_i$ . As such, this naturally translates to an adversarial training paradigm where the filter function and target discriminator are mutual adversaries to each other.

To facilitate the optimization of GetFair, we put forward an alternating training paradigm with Algorithm 1. Specifically, in the first main loop (lines 5-10), mini-batch gradient descent is performed for the target discriminator, which is trained to infer the target labels from the given post embeddings. In the second main loop (lines 11-16), with the discriminator's parameters frozen, the filter function, along with the HSD classifier, are jointly trained with a synergic loss (line 16) that aims to magnify the target classification error while minimizing other loss terms. It is worth noting that, the pretrained transformer-based encoder  $g(\cdot)$  for generating  $s$  will also be finetuned throughout this adversarial training procedure.

## 4 Experiments

To evaluate the efficacy of GetFair in providing accurate yet fair HSD results, we conduct experiments to answer the following research questions (RQs):

**RQ1:** Can GetFair generalize to unseen target groups and outperform state-of-the-art baselines?

**RQ2:** What is the contribution from each core component of it?

**RQ3:** How sensitive GetFair is to its key hyperparameters?

**RQ4:** Is GetFair compatible to different pretrained text encoders?

**Table 1: Hate speech detection and target-aware fairness results with unseen targets. Numbers in bold face are the best results for corresponding metrics, and the second-best results are underlined. Note that we use “↑” and “↓” to indicate higher-is-better and lower-is-better metrics, respectively.**

Dataset	Method	Setting 1						Setting 2					
		Effectiveness (↑)			Fairness (↓)			Effectiveness (↑)			Fairness (↓)		
		Accuracy	F1	AUC	nFPED	nFNED	HF	Accuracy	F1	AUC	nFPED	nFNED	HF
Jigsaw	THSD [50]	0.5947	0.3581	0.5930	0.0074	0.0537	0.0129	0.6591	0.5137	0.6581	<u>0.0013</u>	0.0838	<u>0.0026</u>
	FairReprogram [70]	<u>0.6889</u>	0.5961	<u>0.8040</u>	0.0176	0.0563	0.0268	<u>0.7098</u>	<u>0.6157</u>	<u>0.7098</u>	0.0046	0.0878	0.0087
	SAM [44]	<u>0.6385</u>	<u>0.6850</u>	0.5971	0.1300	0.1354	0.1326	0.6167	0.5148	0.6160	0.1614	0.1538	0.1575
	LWBC [25]	0.5069	<u>0.6446</u>	0.5088	0.0138	<u>0.0207</u>	0.0166	0.4348	0.5231	0.4354	0.0202	<u>0.0671</u>	0.0311
	FEAG [2]	0.6235	0.4162	0.6358	<b>0.0027</b>	<u>0.0358</u>	<u>0.0051</u>	0.6431	0.4591	0.6421	0.0060	0.0861	0.0112
	<b>GetFair</b>	<b>0.7367</b>	<b>0.7262</b>	<b>0.8141</b>	<u>0.0028</u>	<b>0.0087</b>	<b>0.0042</b>	<b>0.7777</b>	<b>0.7719</b>	<b>0.8598</b>	<b>0.0011</b>	<b>0.0569</b>	<b>0.0023</b>
MHS	THSD [50]	0.6315	0.5120	0.6321	<u>0.0071</u>	0.0120	0.0089	0.6422	0.5575	0.6406	0.0340	0.0515	0.0410
	FairReprogram [70]	0.6472	0.5220	<u>0.8134</u>	0.0401	0.0978	0.0569	<u>0.7056</u>	<u>0.7123</u>	<b>0.7780</b>	<u>0.0245</u>	<u>0.0118</u>	<u>0.0159</u>
	SAM [44]	<u>0.7052</u>	0.6581	0.6540	0.1727	0.2326	0.1982	0.6351	0.4642	0.6334	0.1309	0.1406	0.1356
	LWBC [25]	0.5049	<b>0.6688</b>	0.6918	<b>0.0064</b>	<b>0.0036</b>	<b>0.0046</b>	0.5957	0.6424	0.5969	0.1079	0.0996	0.1036
	FEAG [2]	0.5923	0.3671	0.6416	0.0201	0.0530	0.0292	0.6691	0.6170	0.4223	0.0511	0.0399	0.0448
	<b>GetFair</b>	<b>0.7112</b>	<u>0.6592</u>	<b>0.8256</b>	0.0081	<u>0.0041</u>	<u>0.0055</u>	<b>0.7066</b>	<b>0.7302</b>	<u>0.7729</u>	<b>0.0019</b>	<b>0.0124</b>	<b>0.0033</b>

## 4.1 Evaluation Datasets

In our experiments, two publicly acquired benchmarks are in use, namely Jigsaw and MHS. To center our test around generalizability to unseen targets, we hold out two targeted groups for evaluation and use the rest for training and validation. For a thorough comparison, we have created *two test settings* on both datasets by choosing different seen and unseen targets. In what follows, we provide a brief overview of both datasets and their target settings below.

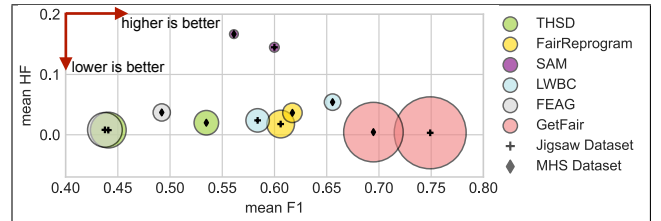
**Jigsaw.** This dataset was released by Google’s online content moderation arm, namely Jigsaw<sup>3</sup> to encourage relevant research. A subset of the posts have been tagged with a variety of identity attributes (i.e., targets) based on the identities mentioned in each post. For our evaluation, we collate all posts with target tags. The target allocation for two test settings are as follows:

- **Setting 1:** Seen targets (training and validation) – {*Male, Female, Homosexual, Christian, Jewish, Black, Mental Illness*}; Unseen targets (test) – {*Muslim, White*}
- **Setting 2:** Seen targets – {*Male, Homosexual, Christian, Jewish, Mental Illness, Muslim, White*}; Unseen targets – {*Female, Black*}

**MHS.** The measuring hate speech (MHS) corpus was created and released by [22]. MHS combines social media posts from YouTube, Twitter, and Reddit that are manually labelled by crowdsourcing workers from Amazon Mechanical Turk (AMT). The AMT annotators have also marked the targets associated with each post. We also test GetFair’s generalization to unseen targets with two settings:

- **Setting 1:** Seen targets – {*Race, Origin, Age, Gender, Disability*}; Unseen targets – {*Religion, Sexuality*}
- **Setting 2:** Seen targets – {*Race, Origin, Gender, Religion, Sexuality*}; Unseen targets – {*Age, Disability*}

Given the space limit, the main statistics and other preprocessing steps are listed in Appendix A.



**Figure 2: Overall performance visualization with both effectiveness and fairness considered. For each dataset, the mean performance is the average of both settings. The size of each scattered point is proportional to  $\frac{\text{mean F1}}{\text{mean HF}}$ .**

## 4.2 Baselines and Metrics

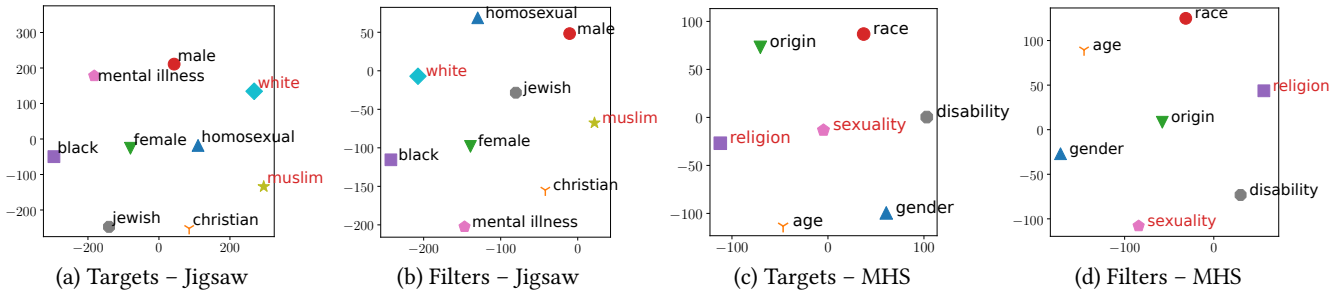
To test the effectiveness of GetFair, we compare it with five fairness-aware baseline methods that aim to debias HSD classifiers, namely THSD [50], FairReprogram [70], SAM [44], LWBC [25], and FEAG [2]. We defer their details to Appendix B.

**Metrics.** For evaluation, we cross-compare all methods from two perspectives: (1) the HSD effectiveness, measured by classification metrics accuracy, F1, and Area under the ROC Curve (AUC); (2) the target-specific fairness, measured by nFPED, nFNED, and HF as in Eq.(1) and Eq.(2).

**Implementation Notes.** For fairness, we adopt the same pre-trained encoder, BERT<sup>4</sup> [23] across all the methods tested for embedding the texts in each post, which is a popular and performant choice in the HSD literature [3, 18]. For GetFair, by default we set discriminator loss coefficient  $\lambda = 0.9$  and the imitation loss weight  $\gamma = 3$  on both datasets according to the contributions and magnitudes of the loss terms. Also, we set  $\mu$  to 0.9 and 0.5, and  $K$  to 1 and 5 respectively for Jigsaw and MHS.  $N$  and  $N'$  are set to 1 and 5 for optimizing GetFair, respectively. We use hidden dimension  $d = 256$  on both datasets. The filter depth  $L$  is consistently set to 1 for optimal efficiency as we do not notice a significant performance

<sup>3</sup><https://www.kaggle.com/competitions/jigsaw-unintended-bias-in-toxicity-classification/data>

<sup>4</sup>We adopt the base version from [https://huggingface.co/transformers/v3.1.0/model\\_doc/bert.html](https://huggingface.co/transformers/v3.1.0/model_doc/bert.html).



**Figure 3: The t-SNE [61] visualization of target indicators and the generated target-specific filter parameters. Targets that are only seen in the test set are annotated in red.**

gain from adding more layers. As for both the target discriminator and HSD classifier, we use the same 3-layer architecture with a consistent 256 hidden dimension in their MLPs.

### 4.3 Overall Performance (RQ1)

We start our analysis with the overall performance. Table 1 reveals the results from all debiasing methods on both datasets. In what follows, we present our observations and findings.

**Effectiveness in HSD.** Table 1 has shown the consistently advantageous HSD effectiveness of GetFair, which has achieved the best accuracy, F1, and AUC on Jigsaw, and is either the best or second-best across these three metrics on MHS. In general, the majority of debiased HSD methods achieve the fairness objective by making a considerable sacrifice in the HSD efficacy, witnessed by the subpar accuracy of LWBC on both Jigsaw and MHS, as well as the low F1 of several models, e.g., THSD and FEAG on both Jigsaw and MHS, as well as FairReprogram on MHS. In contrast, on top of its capability of debiasing, GetFair maintains a high level of utility.

**Fairness in HSD.** On Jigsaw, GetFair yields highly advantageous fairness scores under both settings. Such superiority can also be observed on MHS with Setting 2. On MHS with Setting 1, GetFair has scored the second-best results for both nFNED and HF, only falling behind LWBC by a small margin. In the meantime, it is worth mentioning that the strong fairness of LWBC on MHS is achieved by giving up its utility – its 50% accuracy implies the detection effectiveness is only slightly better than a random classifier. It is also noticed that, among the two metrics nFPED and nFNED, some methods are able to obtain a relatively low score on one of them, but tend to get a substantially higher score on the other. Examples include THSD and FEAG on Jigsaw, where the nFNED scores are one magnitude higher than their nFPED scores, translating into a larger number of real hateful posts being missed out by the trained classifier. Meanwhile, GetFair keeps both nFPED and nFNED low, showcasing a balanced performance. More importantly, considering that the test set contains two unseen targets, this reflects the superior generalizability of GetFair in real-world settings.

**Performance Summary.** To showcase the effectiveness-fairness trade-off of all methods tested, we visualize their performance on both datasets via a scattered plot in Figure 2. Based on the visualization, GetFair is in the highest quartile among all the HSD debiasing methods. Again, this verifies that, compared with other baselines, GetFair is able to achieve the state-of-the-art target-aware fairness results without hurting the real-world usability of the trained HSD classifier. As an additional note on efficiency, with a batch size of

**Table 2: Ablation test with different model architectures.**

Dataset	Architecture	F1	HF
Jigsaw	Default	0.7262	0.0042
	Remove Imitation Learning	0.7189	0.0193
	Remove Semantic Gap Alignment	0.7177	0.0432
	Combinatorial Embedding Filters	0.3853	0.0284
MHS	Default	0.6592	0.0055
	Remove Imitation Learning	0.6439	0.0051
	Remove Semantic Gap Alignment	0.6706	0.0152
	Combinatorial Embedding Filters	0.4866	0.0077

128 and a single Nvidia A40 GPU, GetFair consumes 0.1433s and 0.0962s inference time respectively on Jigsaw and MHS, which can fully support real-time detection.

### 4.4 Ablation Study (RQ2)

We hereby answer RQ2 via ablation study, where we build variants of GetFair by removing/modifying one key component at a time, and the new results from the variants are recorded in Table 2. We use F1 and HF for performance demonstration, and test with Setting 1 on both datasets. Specifically, we are interested in the contributions of imitation learning, semantic gap alignment, and the design of multi-target filters to both the HSD effectiveness and fairness. In what follows, we introduce the corresponding variants and analyze their performance implications.

**Remove Imitation Learning.** The imitation learning objective defined in Eq.(12) aims to uplift the HSD classifier’s performance by aligning the predicted probability distributions generated from the same post’s filtered and unfiltered embeddings. After removing this component, GetFair has experienced a slight drop in F1 scores on both datasets, which showcases the efficacy of imitation learning for improving the detection accuracy. It is also noticed that while the HF score stays stable on MHS, there is an increase in HF on the Jigsaw dataset. One possible reason is that, the performance decrease has also amplified the model’s tendency of producing false negative and false positive results, thus bumping up the HF score.

**Remove Semantic Gap Alignment.** With the semantic gap alignment regularizer (Eq.(7)) removed, the fairness of the HSD results has significantly deteriorated. As the regularization essentially rectifies the way a target-specific filter’s parameters are generated, its removal has incurred a lower quality of generated filters and consequently inferior generalizability to unseen targets. On MHS, a small improvement of F1 is observed, which may attribute to

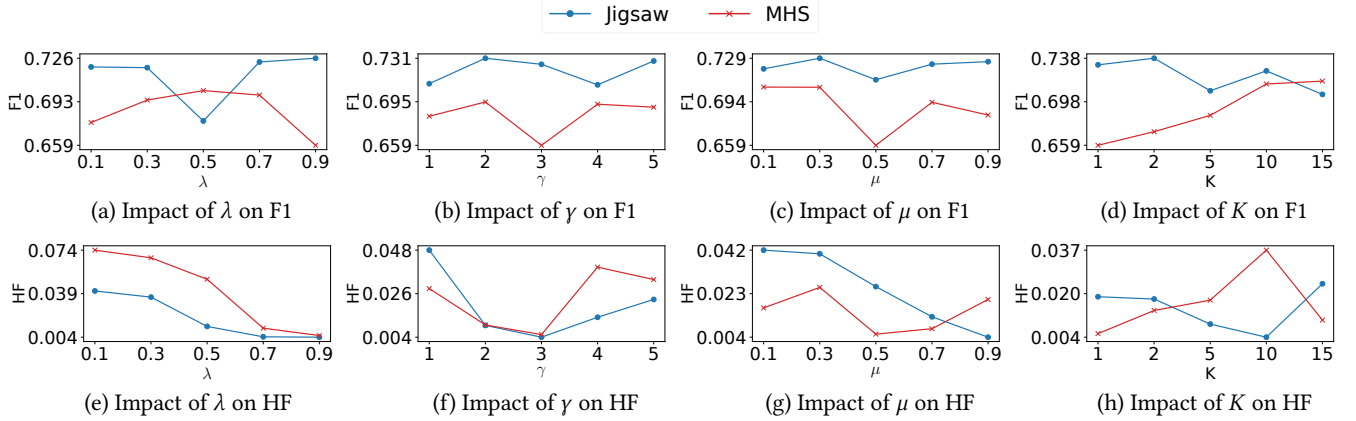


Figure 4: Analysis of the impact from key hyperparameters, with effectiveness and fairness metrics F1 and HF, respectively.

the use of target-specific features as a result of less effective filtering. Besides, to further showcase the efficacy of the semantic gap alignment, we have visualized both the target indicators and their corresponding filter weights generated by the hypernetwork (with the semantic gap alignment in place) by projecting them onto a two-dimensional space via t-SNE [61] via Figure 3. As can be seen, the distances among the target indicators (i.e., raw inputs of the hypernetwork) are mostly kept in their corresponding filter parameters (i.e., outputs of the hypernetwork), thus explaining the fairness boost from this semantic gap alignment scheme.

**Replacing Parameter Ensemble with Combinatorial Embedding Filters.** As we put forward a filter design explicitly for the co-existence of multiple targets within one post, we compare its performance with a straightforward counterpart, i.e., simply performing sum pooling on each individually filtered post embeddings as described in Section 3.1. That is to say, Eq.(5) is updated to  $\tilde{s}_i = \frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} \tilde{s}_{i,t}$ . With this change, a dramatic drop in both the HSD accuracy and fairness is observed. This has verified that, compared with combining outputs from all single-target filters, our parameter ensemble design in GetFair is more capable of preventing the filtered post embedding from noises.

#### 4.5 Hyperparameter Sensitivity (RQ3)

In this section, we examine the impact of some core hyperparameters, namely the discriminator weight  $\lambda$ , imitation weight  $\gamma$ , and regularization weight  $\mu$  used by the synergic loss in Algorithm 1, as well as the rank  $K$  (Section 3.1). This is done by adjusting one hyperparameter at a time and recording the new results achieved, while all other hyperparameters are kept to the default setting. This part of the experiments are also conducted with Setting 1 on both datasets, and similar trends can be observed with Setting 2.

**Impact of  $\lambda$ .** The coefficient  $\lambda$  is tuned in  $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ . As it essentially links to the optimization of target-specific filtering,  $\lambda$  has a stronger impact on HF than on F1, where HF benefits from a larger value of  $\lambda$ . When  $\lambda$  increases from 0.1 to 0.5, there is a rapid improvement in the classification fairness as per the HF score. When  $\lambda$  is sufficiently large (i.e.,  $\lambda \geq 0.7$  in our case), the fairness gain remains positive but appears to be at a lower rate.

**Impact of  $\gamma$ .** Interestingly, the variation of imitation loss weight  $\gamma$  not only causes fluctuations in the HSD effectiveness, but also

Table 3: Performance of GetFair when paired with different pretrained text encoders.

Dataset	Pretrained Encoder $g(\cdot)$	F1	HF
Jigsaw	BERT-base (Default)	0.7262	0.0042
	DistilGPT2	0.6988	0.0051
	RoBERTa-base	0.7501	0.0071
MHS	BERT-base (Default)	0.6592	0.0055
	DistilGPT2	0.6879	0.0055
	RoBERTa-base	0.7864	0.0034

correlates to different fairness levels. As GetFair primarily puts more value on fairness in the HSD,  $\gamma = 3$  is a reasonable choice for a balanced fairness-accuracy trade-off.

**Impact of  $\mu$ .** The regularization weight controls the diversity and quality of the hypernetwork-generated target filters, thus higher sensitivity to  $\mu$  is observed on the fairness metric HF. The general trend is that, as  $\mu$  grows, fairer detection outcomes are attained. When taking the corresponding F1 into account,  $\mu = 0.9$  and  $\mu = 0.5$  are respectively the most sensible settings for Jigsaw and MHS.

**Impact of  $K$ .** The low-rank parameterization of the filter parameters generated by the hypernetwork intends to lower the memory cost while ensuring an adequate level of efficacy. Intuitively, a larger  $K$  is able to preserve more expressiveness for the generated filters. On Jigsaw, a larger  $K$  generally contributes to lower F1 but better HF, while the trend on MHS is on the opposite side. A possible reason is that, on Jigsaw, a stronger filter function removes more target-related information from the post embedding, but also blocks useful features for HSD classification; while on MHS, the spurious target-related features are more implicit, hence a more capable filter function can effectively filter out those noises from the post embedding to achieve a higher F1 score but also more false positive/negative predictions (i.e., a higher HF).

#### 4.6 Compatibility with Other Encoders (RQ4)

GetFair is designed to be compatible with different pretrained text encoders  $g(\cdot)$  as backbones. To verify this compatibility, we have tested GetFair by using two other popular pretrained language models (PLMs), namely DistilGPT2 [20] and RoBERTa (base version) [31] as  $g(\cdot)$ . Similar to RQ2 and RQ3, we test on both datasets under Setting 1 and report F1 and HF metrics.



From the results in Table 3, the first conclusion drawn is that GetFair is able to maintain its performance in HSD tasks when it is paired with different pretrained encoders, especially the classification fairness measured by HF. Secondly, as per F1, RoBERTa yields the highest HSD effectiveness among the three backbones, while DistilGPT2 shows limited effectiveness gain compared with the default BERT used in GetFair. We hypothesize that this is attributed to different PLMs' model capacity. As per the PLMs tested, DistilGPT2 has a lower capacity than the BERT-base we have used (89 million and 110 million parameters respectively), hence producing a similar or even lower accuracy than BERT. In the meantime, the better accuracy of RoBERTa aligns with its higher model size (125 million parameters in the base version) and capacity.

## 5 Related Work

In this section, we review the recent advances in fields that are relevant to our work, specifically hate speech detection (HSD) and the fairness aspects of HSD tasks.

### 5.1 Hate Speech Detection

Mining user-generated web content [37, 56, 57] is a long-lasting research area with versatile applications [9, 38, 71], where hate speech detection is one of the most representative lines of work. Hate speech on social media is commonly defined as a language that attacks or diminishes, that incites violence or hate against groups, based on specific characteristics such as physical appearance, religion, gender identity or other [14]. In summary, extracting representative linguistic features from the post texts lies at the core of various HSD methods. Early practices in HSD involve the use of vocabulary dictionaries like Ortony Lexicon [11] to pinpoint potentially hateful keywords, which then evolves to the use of more sophisticated linguistic features like the term frequency with inverse document frequency (TF-IDF) [11],  $N$ -grams [16], and sentiment [39]. Those text-based features can be easily incorporated with the downstream classifiers for hate speech classification. In the most recent line of work, there has been an adoption of more complex features like images [3, 24] and social connections [34, 46] for the HSD task. Meanwhile, with the rise of language models, especially variations of the transformer family [43, 47, 62], features extracted from the raw texts (a.k.a. embeddings) by those pretrained language models are arguably the most widely adopted option in HSD [3, 18, 24, 28]. The commonly shared goal of HSD is usually performance-oriented, where many new aspects in HSD are attracting an increasing amount of attention, such as efficiency and timeliness of the predictions [60], explainability and transparency of the classification mechanism [35, 65], and robustness of the HSD classifier to adversarial attacks or low-quality data [53]. In what follows, we discuss an emerging research topic in this area, i.e., ensuring the fairness in HSD.

### 5.2 Fairness in Hate Speech Detection

A series of research has uncovered that methods in natural language processing (NLP) tasks are subject to a variety of fairness issues [32, 40, 44], and HSD is no exception as a typical NLP task. Technically, in the context of HSD, the biases can come from either the source [1, 6] (e.g., authors, annotators, and data collectors) and the target [44,

50] of online posts. In this work, our main focus is to demote biases toward the targets of online posts.

The main objective of addressing biases against targeted groups is to diminish the emphasis placed on the inclusion/exclusion of particular terms and redirect attention towards the broader context within the content being evaluated. As discussed in Section 1, data-centric solutions are a representative line of work. Some solutions reweigh each training sample based on its likelihood of introducing bias into the model [48, 69, 72], while some provide additional meta-data, human annotations, or data augmentation [12, 21, 44, 49] for a more rigorous model training process. To bypass the reliance on empiric and human involvement, model-centric solutions aim at removing information related to the spurious target-related features from learned representations of an online post. This is usually achieved by utilizing a dedicated filter module, which is trained end-to-end along with the HSD classifier [10, 18, 19, 26]. It is also seen that some solutions are able to handle intersectional bias between multiple targets [33, 58]. However, the aforementioned methods all suffer from restricted generalizability, as they are trained with the assumption that all targets are seen in the training stage. Due to the distributional discrepancies among different targets, the debiasing effectiveness can hardly transfer to completely unseen targets during inference, creating a practicality bottleneck for real-world applications. Despite the efforts on some generalizable HSD methods [5, 45], their goal is to maintain the HSD accuracy under distributional/domain shift, thus being unable to solve the generalizability challenges associated with target-aware fairness.

## 6 Conclusion

In this paper, to address the deficiency of existing debiasing/fairness-aware HSD methods when handling unseen targets during training, we propose GetFair, which achieves generalizable target-aware fairness in HSD by adaptively generating target-specific filters via a hypernetwork instead of training individualized ones. A suite of innovative designs including low-rank parameterization, semantic gap alignment, and imitation learning are proposed for lowering the memory cost, regularizing the generalizability of target-specific filters, and enhancing the HSD classifier, respectively. Through a series of experiments, we have validated the effectiveness, fairness, and generalizability of GetFair, proving it to be a viable solution to ensuring target-aware fairness in HSD. Possible extensions of GetFair in our future work include discovering time-sensitive patterns [7, 8] for dynamic debiasing, as well as developing lightweight variants [67] of GetFair to further enhance scalability.

## Acknowledgement

This work is supported by Australian Research Council under the streams of Discovery Project (Grant No. DP240101108 and DP240101814), Future Fellowship (No. FT210100624), Discovery Early Career Researcher Award (No. DE230101033 and DE220101597), Center of Excellence (No. CE200100025) and Industrial Transformation Training Centre (No. IC200100022). This work is partially supported by the Swiss National Science Foundation (Contract No. CRSII5\_205975). Financial support from The University of Queensland School of Business in the UQBS 2023 Research Project is gratefully acknowledged.

## References

- [1] Pinkesh Badjatiya, Manish Gupta, and Vasudeva Varma. 2019. Stereotypical bias removal for hate speech detection task using knowledge-based generalizations. In *WWW*. 49–59.
- [2] Parikshit Bansal and Amit Sharma. 2023. Controlling Learned Effects to Reduce Spurious Correlations in Text Classifiers. *ACL* (2023), 2271–2287.
- [3] Aashish Bhandari, Siddhant B Shah, Surendrabikram Thapa, Usman Naseem, and Mehwish Nasim. 2023. CrisisHateMM: Multimodal Analysis of Directed and Undirected Hate Speech in Text-Embedded Images From Russia-Ukraine Conflict. In *CVPR*. 1993–2002.
- [4] Kush Bhatia, Himanshu Jain, Purushottam Kar, Manik Varma, and Prateek Jain. 2015. Sparse local embeddings for extreme multi-label classification. *NeurIPS* 28 (2015).
- [5] Tom Bourgeade, Patricia Chiril, Farah Benamara, and Véronique Moriceau. 2023. What Did You Learn To Hate? A Topic-Oriented Analysis of Generalization in Hate Speech Detection. In *EACL*. 3477–3490.
- [6] Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 6334 (2017), 183–186.
- [7] Tong Chen, Hongzhi Yin, Hongxu Chen, Lin Wu, Hao Wang, Xiaofang Zhou, and Xue Li. 2018. TADA: trend alignment with dual-attention multi-task recurrent neural networks for sales prediction. In *ICDM*. 49–58.
- [8] Tong Chen, Hongzhi Yin, Quoc Viet Hung Nguyen, Wen-Chih Peng, Xue Li, and Xiaofang Zhou. 2020. Sequence-Aware Factorization Machines for Temporal Predictive Analytics. *ICDE* (2020).
- [9] Tong Chen, Hongzhi Yin, Guanhua Ye, Zi Huang, Yang Wang, and Meng Wang. 2020. Try this instead: Personalized and interpretable substitute recommendation. In *SIGIR*. 891–900.
- [10] Lu Cheng, Ahmadreza Mosallanezhad, Yasin N Silva, Deborah L Hall, and Huan Liu. 2022. Bias mitigation for toxicity detection via sequential decisions. In *SIGIR*. 1750–1760.
- [11] Karthik Dinakar, Roi Reichart, and Henry Lieberman. 2011. Modeling the detection of textual cyberbullying. In *AAAI Conference on Web and Social Media*, Vol. 5. 11–17.
- [12] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *AIES*. 67–73.
- [13] Mai ElSherief, Vivek Kulkarni, Dana Nguyen, William Yang Wang, and Elizabeth Belding. 2018. Hate lingo: A target-based linguistic analysis of hate speech in social media. In *ICWSM*, Vol. 12.
- [14] Paula Fortuna and Sérgio Nunes. 2018. A survey on automatic detection of hate speech in text. *Comput. Surveys* 51, 4 (2018), 1–30.
- [15] Tanmay Garg, Sarah Masud, Tharun Suresh, and Tanmoy Chakraborty. 2023. Handling bias in toxic speech detection: A survey. *Comput. Surveys* 55, 13s (2023), 1–32.
- [16] Edel Greevy and Alan F Smeaton. 2004. Classifying racist texts using a support vector machine. In *SIGIR*. 468–469.
- [17] Chuan Guo, Ali Mousavi, Xiang Wu, Daniel N Holtmann-Rice, Satyen Kale, Sashank Reddi, and Sanjiv Kumar. 2019. Breaking the glass ceiling for embedding-based classifiers for large output spaces. *NeurIPS* 32 (2019).
- [18] Soumyajit Gupta, Sooyong Lee, Maria De-Arteaga, and Matthew Lease. 2023. Same Same, But Different: Conditional Multi-Task Learning for Demographic-Specific Toxicity Detection. In *WWW*. 3689–3700.
- [19] Lukas Hauzenberger, Shahed Masoudian, Deepak Kumar, Markus Schedl, and Navid Rekasaz. 2023. Modular and on-demand bias mitigation with attribute-removal subnetworks. In *Findings of the ACL*. 6192–6214.
- [20] HuggingFace. 2019. DistilGPT2. <https://huggingface.co/distilgpt2> (2019).
- [21] Brendan Kennedy, Xisen Jin, Aida Mostafazadeh Davani, Morteza Dehghani, and Xiang Ren. 2020. Contextualizing Hate Speech Classifiers with Post-hoc Explanation. In *ACL*. 5435–5442.
- [22] Chris J Kennedy, Geoff Bacon, Alexander Sahn, and Claudia von Vacano. 2020. Constructing interval variables via faceted Rasch measurement and multitask deep learning: a hate speech application. *arXiv preprint arXiv:2009.10277* (2020).
- [23] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*. 4171–4186.
- [24] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes. *NeurIPS* 33 (2020), 2611–2624.
- [25] Nayeong Kim, Sehyun Hwang, Sungsoo Ahn, Jaesik Park, and Suha Kwak. 2022. Learning debiased classifier with biased committee. *NeurIPS* 35 (2022), 18403–18415.
- [26] Deepak Kumar, Oleg Lesota, George Zerveas, Daniel Cohen, Carsten Eickhoff, Markus Schedl, and Navid Rekasaz. 2023. Parameter-efficient Modularised Bias Mitigation via AdapterFusion. In *EACL*. 2730–2743.
- [27] John P Lalor, Yi Yang, Kendall Smith, Nicole Forsgren, and Ahmed Abbasi. 2022. Benchmarking intersectional bias in NLP. In *NAACL*. 3598–3609.
- [28] Roy Ka-Wei Lee, Rui Cao, Ziqing Fan, Jing Jiang, and Wen-Haw Chong. 2021. Disentangling hate in online memes. In *ACM MM*. 5138–5147.
- [29] Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. 2021. Towards personalized fairness based on causal notion. In *SIGIR*. 1054–1063.
- [30] Zachary Lipton, Julian McAuley, and Alexandra Chouldechova. 2018. Does mitigating ML’s impact disparity require treatment disparity? *NeurIPS* 31 (2018).
- [31] Yinhan Liu et al. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [32] Karthic Madanagopal and James Caverlee. 2023. Bias Neutralization in Non-Parallel Texts: A Cyclic Approach with Auxiliary Guidance. In *EMNLP*. 14265–14278.
- [33] Gaurav Maheshwari, Aurélien Bellet, Pascal Denis, and Mikaela Keller. 2023. Fair Without Leveling Down: A New Intersectional Fairness Definition. In *EMNLP*.
- [34] Sarah Masud, Subhabrata Dutta, Sakshi Makkar, Chhavi Jain, Vikram Goyal, Amitava Das, and Tanmoy Chakraborty. 2021. Hate is the new infodemic: A topic-aware modeling of hate speech diffusion on twitter. In *ICDE*. 504–515.
- [35] Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. Hatexplain: A benchmark dataset for explainable hate speech detection. In *AAAI*, Vol. 35. 14867–14875.
- [36] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM computing surveys* 54, 6 (2021), 1–35.
- [37] Quoc Viet Hung Nguyen, Chi Thang Duong, Thanh Tam Nguyen, Matthias Weidlich, Karl Aberer, Hongzhi Yin, and Xiaofang Zhou. 2017. Argument discovery via crowdsourcing. *VLDBJ* 26 (2017), 511–535.
- [38] Thanh Tam Nguyen, Chi Thang Duong, Matthias Weidlich, Hongzhi Yin, and Quoc Viet Hung Nguyen. 2017. Retaining data from streams of social platforms with minimal regret. In *IJCAI*.
- [39] Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. 2016. Abusive language detection in online user content. In *WWW*. 145–153.
- [40] Ali Omrani, Alireza Salkhordeh Ziabari, Charles Yu, Preni Golazizian, Brendan Kennedy, Mohammad Atari, Heng Ji, and Morteza Dehghani. 2023. Social-group-agnostic bias mitigation via the stereotype content model. In *ACL*.
- [41] Nedjma Ousidhoum, Zizheng Lin, Hongming Zhang, Yangqiu Song, and Dit-Yan Yeung. 2019. Multilingual and Multi-Aspect Hate Speech Analysis. In *EMNLP*. 4675–4684.
- [42] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *EMNLP*. 1532–1543.
- [43] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving Language Understanding by Generative Pre-Training. (2018).
- [44] Alan Ramponi and Sara Tonelli. 2022. Features or spurious artifacts? Data-centric baselines for fair and robust hate speech detection. In *NAACL*. 3027–3040.
- [45] Tharindu Ranasinghe and Marcos Zampieri. 2020. Multilingual Offensive Language Identification with Cross-lingual Embeddings. In *EMNLP*. 5838–5844.
- [46] Manoel Ribeiro, Pedro Calais, Yuri Santos, Virgílio Almeida, and Wagner Meira Jr. 2018. Characterizing and detecting hateful users on twitter. In *AAAI Conference on Web and Social Media*, Vol. 12.
- [47] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- [48] Tal Schuster, Darsh Shah, Yun Jie Serene Yeo, Daniel Roberto Filizzola Ortiz, Enrico Santus, and Regina Barzilay. 2019. Towards Debiasing Fact Verification Models. In *EMNLP*. 3419–3425.
- [49] Indira Sen, Mattia Samory, Claudia Wagner, and Isabelle Augenstein. 2022. Counterfactually Augmented Data and Unintended Bias: The Case of Sexism and Hate Speech Detection. In *NAACL*. 4716–4726.
- [50] Darsh J Shah, Sinong Wang, Han Fang, Hao Ma, and Luke Zettlemoyer. 2021. Reducing target group bias in hate speech detectors. *arXiv preprint arXiv:2112.03858* (2021).
- [51] Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. 2021. Personalized federated learning using hypernetworks. In *ICML*. 9489–9502.
- [52] Shivam Sharma, Atharva Kulkarni, Tharun Suresh, Himanshi Mathur, Preslav Nakov, Md Shad Akhtar, and Tanmoy Chakraborty. 2023. Characterizing the Entities in Harmful Memes: Who is the Hero, the Villain, the Victim?. In *EACL*. 2149–2163.
- [53] Paaras Sheth, Tharindu Kumarage, Raha Moraffah, Aman Chadha, and Huan Liu. 2023. Peace: Cross-platform hate speech detection—a causality-guided framework. In *ECML-PKDD*. 559–575.
- [54] Paaras Sheth, Raha Moraffah, Tharindu S Kumarage, Aman Chadha, and Huan Liu. 2024. Causality guided disentanglement for cross-platform hate speech detection. In *WSDM*. 626–635.
- [55] Leandro Silva, Mainack Mondal, Denzil Correa, Fabricio Benevenuto, and Ingmar Weber. 2016. Analyzing the targets of hate in online social media. In *ICWSM*, Vol. 10. 687–690.
- [56] Xiangguo Sun, Hongzhi Yin, Bo Liu, Qing Meng, Jiuxin Cao, Alexander Zhou, and Hongxu Chen. 2022. Structure learning via meta-hyperedge for dynamic rumor detection. *TKDE* (2022).
- [57] Nguyen Thanh Tam, Matthias Weidlich, Bolong Zheng, Hongzhi Yin, Nguyen Quoc Viet Hung, and Bela Stantic. 2019. From anomaly detection to rumour detection using data streams of social platforms. *VLDB* 12, 9 (2019), 1016–1029.

- [58] Yi Chern Tan and L Elisa Celis. 2019. Assessing social and intersectional biases in contextualized word representations. *NeurIPS* 32 (2019).
- [59] Serra Sinem Tekiroğlu, Yi-Ling Chung, and Marco Guerini. 2020. Generating Counter Narratives against Online Hate Speech: Data and Strategies. In *ACL*. 1177–1190.
- [60] Thanh Tran, Yifan Hu, Changwei Hu, Kevin Yen, Fei Tan, Kyumin Lee, and Se Rim Park. 2020. HABERTOR: An Efficient and Effective Deep Hatespeech Detector. In *EMNLP*. 7486–7502.
- [61] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *JMLR* 9, 11 (2008).
- [62] Ashish Vaswani et al. 2017. Attention is all you need. In *NIPS*. 5998–6008.
- [63] Johannes von Oswald, Christian Henning, Benjamin F Grewe, and João Sacramento. 2020. Continual learning with hypernetworks. In *ICLR*.
- [64] Bencheng Yan, Pengjie Wang, Kai Zhang, Feng Li, Hongbo Deng, Jian Xu, and Bo Zheng. 2022. Apg: Adaptive parameter generation network for click-through rate prediction. *NeurIPS* 35 (2022), 24740–24752.
- [65] Yongjin Yang, Joonkee Kim, Yujin Kim, Namgyu Ho, James Thorne, and Seyoung Yun. 2023. HARE: Explainable Hate Speech Detection with Step-by-Step Reasoning. In *EMNLP*. *ACL*.
- [66] Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W Cohen. 2018. Breaking the Softmax Bottleneck: A High-Rank RNN Language Model. In *ICLR*.
- [67] Hongzhi Yin, Liang Qu, Tong Chen, Wei Yuan, Ruiqi Zheng, Jing Long, Xin Xia, Yuhui Shi, and Chengqi Zhang. 2024. On-Device Recommender Systems: A Comprehensive Survey. *arXiv preprint arXiv:2401.11441* (2024).
- [68] Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. Predicting the Type and Target of Offensive Posts in Social Media. In *NAACL-HLT*. 1415–1420.
- [69] Guanhua Zhang, Bing Bai, Junqi Zhang, Kun Bai, Conghui Zhu, and Tiejun Zhao. 2020. Demographics Should Not Be the Reason of Toxicity: Mitigating Discrimination in Text Classifications with Instance Weighting. In *ACL*. 4134–4145.
- [70] Guanhua Zhang, Yihua Zhang, Yang Zhang, Wenqi Fan, Qing Li, Sijia Liu, and Shiyu Chang. 2022. Fairness reprogramming. *NeurIPS* 35 (2022), 34347–34362.
- [71] Shijie Zhang, Hongzhi Yin, Tong Chen, Zi Huang, Quoc Viet Hung Nguyen, and Lizhen Cui. 2022. Pipattack: Poisoning federated recommender systems for manipulating item promotion. In *WSDM*. 1415–1423.
- [72] Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin Choi, and Noah Smith. 2021. Challenges in Automated Debiasing for Toxic Language Detection. In *ECAL*.

## Appendix

### A Dataset Statistics

As discussed in Section 4.1, we provide the detailed statistics of our experimental datasets, namely Jigsaw and MHS. Both datasets are inherently imbalanced, where the non-hateful posts outnumber hateful ones. Specifically, we have balanced the binary classes of hateful and non-hateful (neutral) posts in the validation and test sets to ensure a more rigorous evaluation. Note that there can be more than one identified target associated with each post. The posts in the test set are allowed to also contain targets that are seen during training, while the training set does not have any posts that mention the two hold-out targets.

Table 4 lists the key statistics of the Jigsaw dataset after being processed for the two different unseen target settings. Analogously, Table 5 below corresponds to the MHS dataset, also with two different unseen target settings.

### B An Overview of Baselines

Following Section 4.2, we hereby provide an extended summary of the baseline methods tested:

- **THSD**: It stands for the token level hate sense disambiguation (THSD) approach proposed by Meta AI Research [50], which regularizes the HSD via token-level consensus to make the classification rely less on target-specific signals.
- **FairReprogram**: This approach appends to the input a set of learnable perturbations called the fairness trigger [70] to achieve the fairness objective under a min-max formulation.
- **SAM**: It is a data-centric approach that utilizes spurious artifact masking (SAM) [44] on the input tokens when training the HSD towards both the accuracy and fairness goals.
- **LWBC**: An ensemble method named learning with biased committee (LWBC) [25] is proposed to adaptively discover and reweight bias-conflicting samples during training.
- **FEAG**: This is the state-of-the-art debiasing method for text classifiers [2] that leverages feature effect augmentation (FEAG) to counter the spurious correlations learned.

**Table 4: Statistics of the Jigsaw dataset.**

Setting	Split	#Hateful Posts	#Non-hateful Posts	Targets
1	Train	11,825	98,396	<i>Male, Female, Homosexual, Christian, Jewish, Black, Mental Illness</i>
	Validation	3,939	3,826	Same as above
	Test	5,973	6,031	+ <i>Muslim, White</i>
2	Train	11,734	76,844	<i>Male, Jewish Homosexual, Christian, Muslim, White, Mental Illness</i>
	Validation	3,764	4,104	Same as above
	Test	5,573	5,608	+ <i>Female, Black</i>

**Table 5: Statistics of the MHS dataset.**

Setting	Split	#Hateful Posts	#Non-hateful Posts	Targets
1	Train	4,118	14,585	<i>Race, Origin, Age, Gender, Disability</i>
	Validation	1,480	1,474	Same as above
	Test	1,496	1,489	+ <i>Religion, Sexuality</i>
2	Train	6,782	24,610	<i>Race, Origin, Gender, Religion, Sexuality</i>
	Validation	2,566	2,707	Same as above
	Test	277	282	+ <i>Age, Disability</i>