

# Cross-Platform Hate Speech Detection with Weakly Supervised Causal Disentanglement

Paras Sheth<sup>1</sup>, Tharindu Kumarage<sup>\*1</sup>, Raha Moraffah<sup>\*1</sup>, Aman Chadha<sup>\*\*2,3</sup>,  
and Huan Liu<sup>1</sup>

<sup>1</sup> Arizona State University, Tempe, AZ, USA

<sup>2</sup> Stanford University, Stanford, CA, USA

<sup>3</sup> Amazon Alexa AI, Sunnyvale, CA, USA

{psheth5, kskumara, rmoraffa, huanliu}@asu.edu  
hi@aman.ai

**Abstract.** Content moderation faces a challenging task as social media’s ability to spread hate speech contrasts with its role in promoting global connectivity. With rapidly evolving slang and hate speech, the adaptability of conventional deep learning to the fluid landscape of online dialogue remains limited. In response, causality-inspired disentanglement has shown promise by segregating platform-specific peculiarities from universal hate indicators. However, its dependency on available ground-truth target labels for discerning these nuances faces practical hurdles with the incessant evolution of platforms and the mutable nature of hate speech. Using confidence-based reweighting and contrastive regularization, this study presents HATE-WATCH, a novel framework of weakly supervised causal disentanglement that circumvents the need for explicit target labeling and effectively disentangles input features into invariant representations of hate. Empirical validation across platforms (two with target labels and two without) positions HATE-WATCH as a novel method in cross-platform hate speech detection with superior performance. HATE-WATCH advances scalable content moderation techniques towards developing safer online communities.

## 1 Introduction

**Warning:** *This paper contains contents that may be offensive or upsetting.*

Hate speech on social media not only disrupts healthy dialogue but also threatens societal harmony and safety. It has been linked to real-world violence, as seen in events like the 2019 Christchurch mosque shootings in New Zealand, where the perpetrator was influenced by online extremist content [25]. Similarly, the COVID-19 pandemic saw a spike in hate crimes against Asians, fueled by xenophobic comments online [8]. These instances underscore the urgent need for advanced moderation tools to curb the spread of online hate and prevent its escalation into physical aggression and division.

The identification of hate speech on social media poses considerable hurdles because of the rapid expansion of these platforms [14] and the emergence of distinct community standards. Lacking platform-specific labeled data, new platforms frequently

\* Both authors contributed equally.

\*\* Work does not relate to the position at Amazon

rely on models that were trained using data from existing platforms [10] or by directly leveraging large language models (LLMs) [16]. Utilizing these solutions has pitfalls as LLMs might not be able to determine hate due to the social-context demands of the task [44]. On the other hand, the hate speech models trained on one platform may not generalize to other platforms due to various factors, including; (1) the variance in regulatory policies, from stringent content controls on platforms like FoxNews [5] to more lenient policies on sites like GAB [3], and (2) hate speech’s context-dependent characteristics. For example, the COVID-19 pandemic’s increase in hate speech against Asians brought to light the limitations of pre-pandemic models, which lacked the contextually relevant training data to effectively identify and mitigate such emerging forms of hate [13]. These complexities illustrate the challenges in creating adaptable and accurate hate speech detection models across varying social media environments.

Efforts to enhance cross-platform model performance have explored linguistic cues [31] and context-aware fine-tuning [7]. Yet, these methods often fall short, either forming spurious correlations or embedding platform-specific biases that hinder generalization. Recognizing these limitations, recent causality-based approaches have emerged as promising alternatives, leveraging causal cues [35] or graphs [36] to guide models in learning platform-invariant hate representations. While these causality-based methods address several of the aforementioned challenges, they introduce significant demands for data annotation and comprehensive auxiliary data, such as the precise target of hatred, raising concerns regarding feasibility, privacy, and ethical implications.

In response to these challenges and building upon the foundation laid by causality-based approaches, we propose a novel **Weakly supervised cAusality-aware disenTanglement** framework for **Cross-platform Hate speech detection** named HATE-WATCH<sup>4</sup>, that synthesizes weak supervision with contrastive learning to minimize the dependency on extensively labeled target data. Inspired by advancements in fine-tuning language models (LMs) with weak supervision [43], HATE-WATCH aims to mitigate the challenges of noisy labels and sparse data, significantly enhancing the model’s adaptability across various social media platforms. By leveraging contrastive learning HATE-WATCH facilitates distinguishing between different target representations. Furthermore, with confidence-based sample reweighting methods HATE-WATCH suppresses label noise propagation during training ensuring effective model learning. This approach not only simplifies the demands for auxiliary data but also leverages the inherent variability and noise within weakly labeled datasets as an asset, thereby improving the robustness and generalization of hate speech detection.

HATE-WATCH represents a significant advancement in cross-platform hate speech detection research by offering a scalable and equitable solution that addresses the practical and ethical concerns associated with current methodologies. Our contributions are summarized as follows:

- We introduce a causality-based detection framework HATE-WATCH enhanced by weak supervision and contrastive learning, significantly reducing the reliance on labeled data for platform-dependent features.

<sup>4</sup> The code can be accessed from [https://anonymous.4open.science/r/HATE-WATCH-406C/HATE\\_WATCH.ipynb](https://anonymous.4open.science/r/HATE-WATCH-406C/HATE_WATCH.ipynb)

- HATE-WATCH demonstrates superior adaptability and generalization in hate speech detection across four social media environments (two with target labels and two without), validated through comprehensive experiments.

## 2 Related Work

Recent strides in hate speech detection emphasize generalizability across varying social media landscapes and demographics, addressing the challenge posed by the dynamic nature of online hate speech. The diversity in hate speech manifestation complicates detection of emerging topics and groups [29].

**Enhancing Cross-platform Hate Speech Detection via Representation Disentanglement:** Addressing the scarcity of widely available data, recent efforts have harnessed auxiliary features such as user attributes [38] and annotator demographics [42], alongside LMs including Llama 2, GPT-3.5 [19], BERT, and RoBERTA [7,27] for cross-platform hate speech detection. However, the absence of explicit training objectives for generalizability in these detectors could lead to errors in cross-platform scenarios [35]. Recent studies aim to enhance generalizability through explicit modeling of hate representations, using causal cues like aggression and sentiment [35]

Recent advances in representation disentanglement, encoding distinct aspects of data [1], show significant promise in generalizing hate speech detection across modalities by separating hate-target entities [21,41] and distinguishing genuine from spurious content elements [32]. Notably, causality-guided disentanglement [36], inspired by the platform-specific nature of hate targets [26,6], uses invariant features for state-of-the-art generalization. This approach demonstrated that by removing platform-dependent features and leveraging the invariant features one can achieve state-of-the-art generalization, though it often presupposes available hate-target annotations. Therefore, our study aims to explore a more pragmatic approach, where ground-truth target labels are absent.

**Weakly Supervised Learning** In the face of inconsistent data and varying annotations, this strategy refines noisy labels, which is advantageous for hate speech models [15]. The rise of pre-trained LMs has facilitated the emergence of weakly supervised methods, such as Xclass [39] and LOTClass [28], which leverage contextual information from label names and keywords for denoising and learning enhancement. Recent studies have demonstrated the effectiveness of these methods in achieving cross-dataset generalization in hate speech detection [15]. Current developments demonstrate that by leveraging contrastive self-training, a technique that harnesses the denoising capabilities of pre-trained LMs widens their applicability to a broader array of tasks [43]. In line with these advancements, our research adopts the contrastive self-training method to learn cross-platform hate speech representations with weak target labels.

## 3 Methodology

### 3.1 Preliminaries

Separating platform-specific from invariant features is crucial for identifying hate speech across platforms. Causality, a tool proven to boost a model’s generalization [4], can un-

derpin this process, as shown in Figure 1. Here, the input data can be decomposed into (1) a causal component, constant across platforms (e.g., aggression) and (2) a platform-specific component, unique to each platform (e.g., platform policies).

Distinguishing between these factors is pivotal, yet quantifying platform policies is challenging. Previous work suggests using hate targets as proxies for platform-specific aspects [36]. However, acquiring auxiliary and ground-truth labels is difficult, especially on data-scarce platforms. Additionally, trying to disentangle in an unsupervised manner yields incorrect disentanglement rendering the model ineffective [24].

To overcome these barriers, we introduce HATE-WATCH, which bypasses the need for target labels by adopting a weakly supervised, contrastive learning method with confidence-based reweighting for precise disentanglement. Given a source corpus  $D_{source}$  consisting of textual inputs  $X = \{x_1, x_2, \dots, x_n\}$  and corresponding hate speech labels  $Y = \{y_1, y_2, \dots, y_n\}$ , HATE-WATCH aims to establish a mapping  $F : X \rightarrow Y$  that pinpoints the causal attributes of hate speech. It employs contrastive learning to maintain distinct target representations and integrates confidence-based reweighting to reduce noise impact, enhancing model training’s effectiveness.

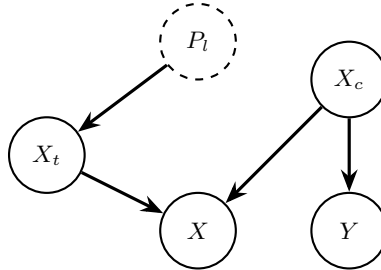


Fig. 1: The causal graph illustrates the hate speech detection mechanism, where  $X_c$  denotes the causal factors predicting hate speech,  $Y$  the hate label,  $X_w$  the target,  $X$  the input, and  $P_l$  the latent platform variable affecting the target.

### 3.2 Disentangling Causal and Target representations

The manifestation of hate speech often intertwines with the targets it seeks to disparage. Even though the precise targets could fluctuate significantly between platforms—a reflection of the platform-dependent nature of hate speech—the fundamental organization and severity of the offensive material frequently show striking consistency. For instance, foul language used to disparage one group can frequently be applied with little modification to another target. This finding forms the basis of our methodology in HATE-WATCH, where we propose that the essential elements of hate—its causal features—can be successfully extracted from the text even in the lack of clear hate target labeling. The goal of HATE-WATCH is to capture the substance of hate speech by abstracting from the changeable, platform-specific targets. This allows to identify and analyze hateful content with an intuitive knowledge that goes beyond specific targets.

HATE-WATCH utilizes a three-component architecture—encoder, disentanglement mechanism, and decoder—mirroring a VAE structure optimized for weak supervision. The encoder  $q_\phi$ , a language model like RoBERTa [23], transforms input text  $x$  into an embedding  $z$  via  $z = q_\phi(\gamma(x))$ , where  $z \in \mathbb{R}^{s_l \times h_d}$  is the embedding output,  $\gamma(x)$  is the tokenization of  $x$ ,  $s_l$  represents sequence length, and  $h_d$  the embedding dimension.

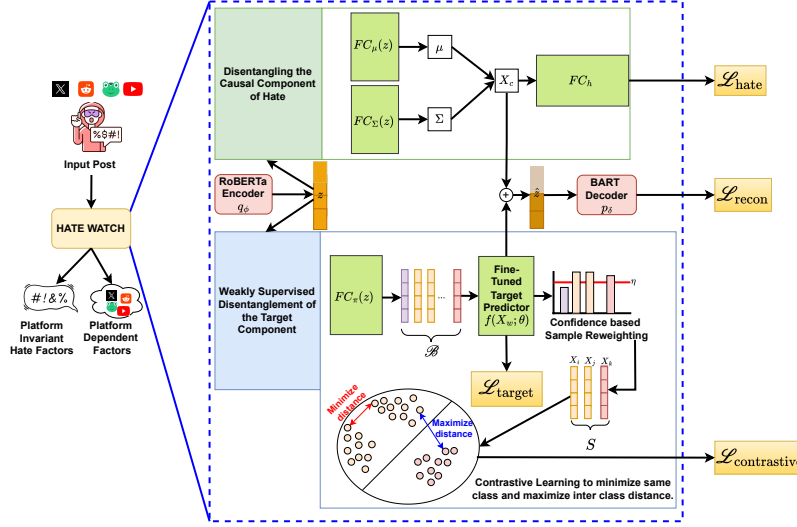


Fig. 2: The HATE-WATCH architecture processes input  $X$  via a RoBERTa to get initial representation  $z$ . This  $z$  undergoes disentanglement into a causal component to identify invariant hate factors  $X_c$ , and a weakly supervised target component  $X_w$  without true target labels. Both components’ outputs,  $X_c$  and  $X_w$ , are merged to form reconstructed embedding  $\hat{z}$ , which is decoded by BART to produce reconstructed input  $\hat{X}$ .

The [CLS] token embedding  $z^{[CLS]} \in \mathbb{R}^{h_d}$ , indicating the input’s overall representation, exemplifies this process. The overall framework of HATE-WATCH is visualized in Figure. 2.

**Disentangling the Causal Component** HATE-WATCH aims to separate the embedding  $z$  in order to extract the subset of causal features  $X_c$  using its VAE design. Disentanglement is a procedure that seeks to derive a factorized representation by aligning individual latent variables with distinct explanatory factors of the data. The separation process is facilitated by the typical Gaussian prior applied to the latent space inside the VAE framework, as well as the posterior approximation achieved by a parameterized neural network. In this study, two feedforward neural networks, namely  $FC_\mu$  and  $FC_\Sigma$ , are utilized to map the variable  $z$  onto the parameters of a Gaussian distribution. This mapping is designed to capture the inherent causal relationship present in hate speech. The distribution is used to create the latent causal feature set  $X_c \in \mathbb{R}^{h_{causal}}$ , employing the re-parameterization trick [17] as follows:

$$X_c = Enc_1(\mu_z, \Sigma_z) = \mu_z + \Sigma_z \odot \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where  $\mu_z = FC_\mu(z)$ ,  $\Sigma_z = FC_\Sigma(z)$ ,  $\epsilon$  is a random noise vector sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , representing the standard normal distribution, and  $\mathbf{I}$  is the unit variance.

**Disentangling the Target Component** In addition to the causal features captured by the embedding  $z$ , our model must contend with platform-dependent features that vary across different social media environments, denoted by  $X_w$ . Given that these fac-

tors are not easily quantifiable (e.g., policies and regulations of a platform) we leverage a proxy, i.e., the target of hate based on the observation that the distribution of the targets is platform dependent [12,40].

Unlike the invariant causal features, these platform-dependent attributes are inherently discrete, with hate speech targeting specific categories such as Race, Religion, Gender, among others [2]. Given the discrete nature of targets and the absence of ground truth labels in various real-world scenarios, we employ a contrastive learning framework along with confidence based weighted resampling, inspired by [43] to learn a distinct representation for each target class.

**Contrastive Learning Framework** The essence of our approach is to enforce a structure within the latent space that brings closer the representations of samples with similar targets while distancing those with different ones. To implement this, we utilize a feed-forward neural network  $FC_\pi$  to map the embedding  $z$  onto a latent space  $X_w$  where,  $X_w = FC_\pi(z)$ , and  $X_w \in \mathbb{R}^{h_{disc}}$  holds the latent representation for the target classes. This representation is further passed through a target classifier  $f(X_w; \theta)$  which is a finetuned feed-forward network, to obtain the soft-labels  $\tilde{y}$ .

Given the soft pseudo-target labels  $\tilde{y}$  generated by the model, we define a contrastive loss that encourages a compact cluster formation for each target class. Specifically, we first select high-confidence samples  $\mathcal{S}$  from  $\mathcal{B}$  as  $\mathcal{S} = \{X_w \in \mathcal{B} \mid w(x) \geq \eta\}$  where  $\eta$  is a predefined threshold. Then, we define the similarity between each pair  $X_{w_i}, X_{w_j} \in \mathcal{S}$  as follows:

$$W_{ij} = \begin{cases} 1, & \text{if } \arg \max_{k \in \mathcal{Y}} [\tilde{y}_i]_k = \arg \max_{k \in \mathcal{Y}} [\tilde{y}_j]_k, \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

where  $\tilde{y}_i, \tilde{y}_j$  are the soft pseudo-labels obtained from for  $X_{w_i}, X_{w_j}$ , respectively. Next, we define the contrastive regularizer as follows:

$$\begin{aligned} \mathcal{L}_{contrastive} &= \sum_{(X_{w_i}, X_{w_j}) \in \mathcal{S} \times \mathcal{S}} [W_{ij} d_{ij}^2 + (1 - W_{ij}) \max(0, \beta - d_{ij})^2], \\ \ell &= W_{ij} d_{ij}^2 + (1 - W_{ij}) \max(0, \beta - d_{ij})^2, \end{aligned} \quad (3)$$

where the contrastive loss  $\mathcal{L}_{contrastive}$  utilizes  $\ell(\cdot, \cdot)$  to measure the Euclidean distance ( $d_{ij}$ ) between  $X_{w_i}, X_{w_j}$ , and  $\beta$  denotes a pre-defined margin. The function  $\ell$  is emphasized for its use in calculating the individual contributions to the overall contrastive loss, as detailed in works such as [9,37].

Eq. 3 penalizes the distance between samples from the same class, i.e.,  $W_{ij} = 1$ ; for samples from different classes, the contrastive loss is considerable if their distance is smaller. The regularizer maintains at least  $\gamma$  distance between dissimilar samples and forces similar samples to be close together in this fashion.

**Confidence-Based Sample Reweighting** The contrastive loss facilitates learning distinct representations for each target class, however in the absence of ground truth labels it may result in label-noise propagation [33]. To mitigate the propagation of label

noise due to the weakly supervised setting, we incorporate a confidence-based sample reweighting mechanism. Each sample  $i$  is assigned a weight based on the entropy of its predicted probability distribution  $\tilde{y}_i$ , reflecting the confidence of the model in its prediction. The weights are calculated as  $w_i = 1 - \frac{H(\tilde{y}_i)}{\log(C)}$ , where  $H(\tilde{y}_i)$  denotes the entropy of the predicted probabilities for sample  $i$ , and  $C$  is the number of target classes.

During contrastive self-training, the sample reweighting strategy encourages high confidence samples. This approach, however, depends on incorrectly classified samples being considered low confidence, which may not be the case unless we stop overly optimistic forecasts. To promote smoothness over forecasts, we utilize a confidence based regularizer as,

$$\mathcal{L}_{conf} = \frac{1}{|\mathcal{S}|} \sum_{X_w \in \mathcal{S}} \mathcal{D}_{KL}(u || f(X_w; \theta)), \quad (4)$$

where  $\mathcal{D}_{KL}$  is the KL-divergence and  $u_i = \frac{1}{C}$  for  $i = 1, 2, \dots, C$ , and  $C$  is the total target classes. Such term constitutes a regularization to prevent over-confident predictions and leads to better generalization [30]. We thus define the loss function as,

$$\mathcal{L}_{target}(\theta, \tilde{y}) = \frac{1}{|\mathcal{S}|} \sum_{X_w \in \mathcal{S}} w(X_w) \mathcal{D}_{KL}(\tilde{y} || f(Z_\pi; \theta)), \quad (5)$$

where  $\mathcal{D}_{KL}$  is the Kullback-Leibler (KL) divergence [18].

By integrating the contrastive strategy with confidence-based sample reweighting, our model HATE-WATCH is capable of learning to disentangle target representations that are robust to variations in data distribution and labeling noise.

**Reconstructing the Input from the Disentangled components** To facilitate the training, we aim to reconstruct the input from the obtained components  $X_c$  and  $X_w$ . To do so, we first concatenate them together as  $[X_c | X_w]$  where  $[|]$  is the concatenation operation. Then we pass it through another feed-forward neural network, namely,  $FC_{\hat{z}}$ . The obtained representation after concatenation is given by  $\hat{z} = FC_{\hat{z}}([X_c | X_w])$ . To recreate the input, we feed  $\hat{z}$  through a LM decoder  $p(x|\hat{z})$ . We utilize BART-base decoder [22] as the LM-decoder, as it shares the vocabulary with RoBERT [23] implementations and is proved powerful in many generative tasks. The obtained tokens are given as  $\hat{x} = LMHead(p_\delta(\hat{z}))$ , where  $LMHead$  is a feed-forward layer to map the decoder  $p_\delta$  embeddings into tokens. The reconstructed loss is computed between the input token ids and the reconstructed token ids and is formulated as follows:

$$L_{recon}(y(x), \hat{x}) = - \sum_{i=1}^{s_l} y(x) \log(\hat{x}_i), \quad (6)$$

where  $s_l$  is the sequence length. We also apply KL-divergence losses for both latent spaces to make sure that the posteriors of the disentangled latent spaces are near to their prior distribution. The disentanglement module’s Evidence Lower Bound (ELBO) is expressed as follows:

$$\mathcal{L}_{VAE} = \mathcal{L}_{recon} + \alpha_t * L_{\mathbb{D}_{target}} + \alpha_c * L_{\mathbb{D}_{causal}}, \quad (7)$$

Dataset	# Posts	Hateful Posts	Hate %	Target?
<b>GAB</b> [27]	11,093	8,379	75.5	✓
<b>Reddit</b> [20]	39,811	15,388	38.6	✗
<b>X</b> [11]	24,802	9,118	36.7	✗
<b>YouTube</b> [34]	1,026	642	62.5	✓

Table 1: Dataset statistics with percentage of hateful comments or posts.

where  $\alpha_t$  represents the coefficient that controls the contribution of the KL loss for the target, and  $\alpha_c$  represents the coefficient that controls the contribution of the causal KL loss. To facilitate learning in absence of ground-truth target labels, and by leveraging confidence based denoising and contrastive learning inspired by [43]  $L_{\mathbb{D}_{target}}$  and  $L_{\mathbb{D}_{causal}}$  are given by,

$$\begin{aligned} L_{\mathbb{D}_{target}} &= \mathcal{L}_{target} + \delta_{cont} \cdot \mathcal{L}_{contrastive} + \delta_{conf} \cdot \mathcal{L}_{conf}, \\ L_{\mathbb{D}_{causal}} &= D_{KL}(Enc_1(X_c | X) || p(X_c)), \end{aligned} \quad (8)$$

### 3.3 Model Training

The disentangled latent causal representation  $X_c$  is used to calculate the classification probability for hate speech detection, given by  $\hat{y}_i = \text{Softmax}(FC_h(X_c))$  where  $FC_h$  represents a fully connected layer for hate classification. The total loss  $\mathcal{L}$  is then calculated using the conventional cross-entropy method:

$$\mathcal{L}_{hate} = -\frac{1}{N} \sum_{i=1}^{|D_{source}|} y_i \log \hat{y}_i \quad (9)$$

where  $D_{source}$  represents the source domain data,  $y_i$  represents the true hate label, and  $\hat{y}_i$  denotes the predicted hate labels. Lastly, we integrate every suggested module and train using a multi-task learning approach:

$$\mathcal{L} = \mathcal{L}_{hate} + \mathcal{L}_{VAE} \quad (10)$$

## 4 Experiments

In this section, we conduct a series of experiments aimed at verifying the capability of HATE-WATCH in acquiring generalizable representations to identify hate speech through causality-aware disentanglement without reliance on target labels. We use benchmark datasets from multiple platforms to guarantee a comprehensive analysis. We aim to answer the following research questions:

- **RQ.1** Is it possible to effectively disentangle the causal and platform dependent factors even in the absence of ground truth target labels?
- **RQ.2** How does forgoing target labels affect disentanglement effectiveness compared to using labels derived from LLMs?
- **RQ.3** Does weakly-supervised disentanglement learn invariant relationships?



#### 4.1 Dataset and Evaluation Metrics

To identify hate speech on English-language benchmark datasets from GAB, Reddit, X, and YouTube, binary classification is used. GAB, which is annotated for hatefulness and is obtained from the GAB website [27]. The YouTube dataset, which has been described in detail by [34], has comments that use offensive language. Hate and target labels are provided by both datasets.

Four ordinal labels are used in Reddit’s hatefulness-categorized dataset [20]. We classify any denigrating information as hateful for consistency’s sake. In a similar vein, content is classified as Hate, Offensive, or Neither in X’s dataset that was generated from tweets [11], with the first two being deemed hateful. X and Reddit are strictly binaryized into hate and non-hate similar to the other platforms. Unlike GAB and YouTube, X and Reddit do not have any hate target labels. These datasets are compiled in Table 1, and their evaluation is based on the macro F1-measure.

#### 4.2 Experimental Settings

**Implementation Details** We trained our framework using RoBERTa-base for the LM-encoder and BART-base for the LM-decoder with the Huggingface Transformers library. We optimized the model with cross-entropy loss and AdamW, using a learning rate of 0.0001, dropout rate of 0.2, and parameters  $\alpha_t$  as 0.05,  $\alpha_c$  as 0.05,  $\delta_{cont}$  as 0.001,  $\delta_{conf}$  as 0.001,  $\eta$  as 0.95 and  $\beta$  as 2. Training was conducted on an NVIDIA GeForce RTX 3090 GPU with 24 GB VRAM, with early-stopping.

**Baselines** In our analysis, the HATE-WATCH framework is evaluated against leading methods across three categories: fine-tuned LMs, causality-aware techniques, and weakly supervised approaches, to benchmark its effectiveness in cross-platform hate speech detection.

**Fine-Tuned LMs: HateXplain** [27] and **HateBERT** [7] are optimized for hate speech identification, the former focusing on a trinary classification with annotated justifications using X and GAB data, and the latter finetuned on 1.5 million Reddit posts for training a BERT-base model.

**Causality-Aware Techniques:** Utilizing sentiment and aggression features (**PEACE** [35]) or auxiliary information like target labels (**CATCH** [36]), these methods aim to improve generalization for hate speech from a causal lens.

**Weakly Supervised Approaches: XClass** [39] and **LOTClass** [28] employ contextual cues and semantic relationships for noise reduction and category prediction, enhancing model training in the absence of explicit labels.

#### 4.3 RQ1. Performance Comparison

Using macro-F1 for assessment, the generalizability of our model, HATE-WATCH is examined via comparative analysis against several baselines and across multi-platform datasets (Table 2). This analysis evaluates the cross-platform performance of HATE-WATCH by emphasizing causal rather than non-causal elements without auxiliary labels. Observations w.r.t RQ.1 are as follows:

Dataset Type	Source	Target	Models							
			HateXplain	HateBERT	PEACE	CATCH	XClass	LOTClass	OURS	
With Target Labels	GAB	GAB	<u>0.87</u>	<b>0.89</b>	0.76	0.82	0.79	0.77	0.81	
		YouTube	0.62	0.6	0.64	<b>0.66</b>	0.55	0.51	<u>0.65</u>	
		X	0.54	0.59	0.6	NA	0.6	0.58	<b>0.63</b>	
	Reddit	Reddit	0.56	<b>0.62</b>	0.61	NA	0.55	0.57	<b>0.62</b>	
		YouTube	GAB	0.47	0.52	0.48	<b>0.56</b>	0.5	0.45	<u>0.53</u>
			YouTube	<b>0.88</b>	0.84	<u>0.86</u>	0.79	0.77	0.72	0.76
	X		0.49	0.53	<u>0.57</u>	NA	0.47	0.45	<b>0.59</b>	
	Reddit	Reddit	0.52	0.54	<u>0.58</u>	NA	0.51	0.49	<b>0.64</b>	
		X	GAB	0.62	0.61	0.6	NA	<u>0.64</u>	0.61	<b>0.65</b>
YouTube			0.61	0.62	<b>0.65</b>	NA	0.6	0.57	<u>0.64</u>	
X	0.93		<u>0.92</u>	<b>0.93</b>	NA	0.81	0.78	0.91		
Reddit	<u>0.51</u>		0.45	0.45	NA	<u>0.51</u>	0.47	<b>0.54</b>		
Without Target Labels	Reddit	GAB	0.53	0.57	<u>0.63</u>	NA	0.55	0.56	<b>0.65</b>	
		YouTube	0.39	0.44	<b>0.56</b>	NA	0.45	0.44	<u>0.55</u>	
		X	<u>0.55</u>	0.49	0.54	NA	0.53	0.54	<b>0.56</b>	
		Reddit	<u>0.89</u>	<b>0.9</b>	<u>0.89</u>	NA	0.79	0.81	0.88	
Average Performance			0.62	0.63	<u>0.65</u>	0.18	0.6	0.58	<b>0.66</b>	

Table 2: Cross-platform and in-dataset evaluation results for the different baseline models compared against HATE-WATCH. Boldfaced values denote the best performance, and the underline denotes the second-best performance. NA implies Not Applicable due to absence of target labels.

- The HATE-WATCH model exhibits notable generalizability across different platforms. Despite the absence of auxiliary labels, e.g. the target of hate labels, HATE-WATCH still manages to discern and prioritize causal over non-causal factors. This indicates that the contrastive learning strategy coupled with the confidence based denoising allows HATE-WATCH to maintain performance in diverse environments, reinforcing the model’s utility in scenarios where hate targets may not be readily available.
- When comparing causal models, **CATCH**, **PEACE**, and HATE-WATCH, we observe that these models, on average, outperform the non-causal methods, highlighting the importance of causality in model generalization. PEACE performs well even in the absence of target labels, however, PEACE utilizes only two causal cues i.e. sentiment and aggression for learning generalizable representations. However, there may be other factors (e.g. user history, demographics) that are influential in determining hatefulness but not easily quantifiable, limiting PEACE’s capabilities. HATE-WATCH, on the other hand, models these cues in the latent space, displaying a competitive edge in settings without target labels, hinting at its efficient use of underlying causal relationships that are consistent across platforms, in contrast to non-causal features that may vary and lead to overfitting.
- As for models like **HateBERT** and **HateXplain**, while powerful and based on robust architectures, they appear to underperform in cross-platform scenarios. This could be indicative of a possible overfitting to platform-specific nuances within the training data, limiting their applicability in a cross-platform context where

Source	Model	GAB	YouTube
GAB	Fully Supervised (CATCH)	<b>0.82</b>	<b>0.66</b>
	Weak-Supervised-GPT4	0.71	0.60
	Unsupervised	0.69	0.59
	HATE-WATCH	<u>0.81</u>	<u>0.65</u>
YouTube	Fully Supervised (CATCH)	<b>0.56</b>	<b>0.79</b>
	Weak-Supervised-GPT4	0.51	<u>0.78</u>
	Unsupervised	0.43	0.55
	HATE-WATCH	<u>0.53</u>	0.76

Table 3: Comparison of HATE-WATCH’s generalization capabilities through macro-F1 against different modeling techniques.

platform-invariant features are essential. Conversely, HATE-WATCH circumvents this by its very design, which is geared towards recognizing and utilizing features of hate speech that are pertinent to hate and shared across various platform.

- Compared to supervised models, both **XClass** and **LOTClass** perform close to supervised models, even in the absence of explicit hate labels indicating how self-training and denoising can facilitate qualitative learning for scenarios where labeled data is scarce.

#### 4.4 RQ2. How Helpful Are The Weak-Labels?

The utilization of weak labels in hate speech detection models serves as an intriguing middle ground between the often resource-intensive fully supervised techniques and the less constrained unsupervised approaches. In this experiment, we investigated the effectiveness of weak labels by contrasting our model—which makes use of these labels—with three different approaches. The first comparison is made with CATCH, a fully supervised technique that allows for a more focused disentanglement by providing both hate and target labels. The second model in our comparative study involves GPT-4, where we crafted a prompt as follows:

##### Prompt to Detect Hate Targets

The following examples show the post and the target group being talked about in the post. Examples: ... Now, given the following posts, identify the main target group of the post. The target category of the post refers to the entity being talked about in the post. The possible categories are Ability/Disability, Class, Gender, Immigration Status, Nationality, Race, Religion, Sexuality, and Sexual Preferences.

In order to bridge the gap between fully supervised and unsupervised algorithms for hate speech identification, weak labels are deployed. We evaluated the effectiveness

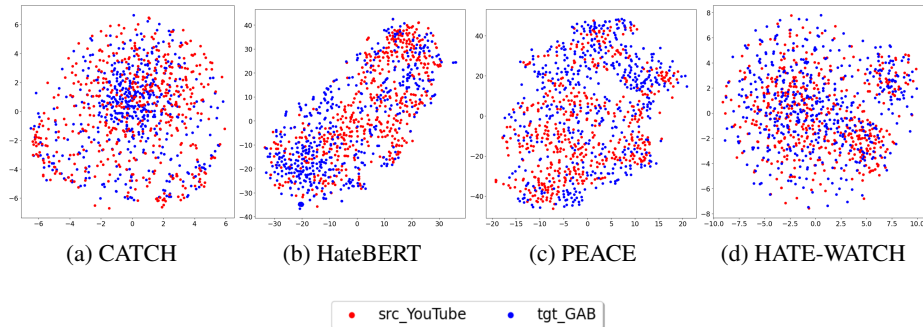


Fig. 3: Visualizing the representations from different models to verify invariance across platforms. `src` (`tgt`) denote the source (target) platforms.

of weak labels using HATE-WATCH in comparison to three approaches: a fully supervised one (CATCH), one that had GPT-4 assistance, and one that was completely unsupervised. As Table 3 shows, GPT-4 performs better than the unsupervised approach, but it falls short of HATE-WATCH.

The training signal may be diluted by noise introduced by the powerful label synthesis of GPT-4, which has a 15-20% target misidentification rate. In comparison, weak label denoising is essential for reliable feature disentanglement in the presence of noisy data in HATE-WATCH. Furthermore, LLMs such as GPT-4 might not be able to properly recognize the contextually complicated nature of hate speech, while HATE-WATCH’s lax monitoring is more effective.

HATE-WATCH outperforms unsupervised and GPT-4-based approaches in terms of generalization within and between the YouTube and GAB domains. Its ability to transfer hate speech characteristics across platforms with remarkable resilience highlights the efficacy of the weak-label framework. The study shows that improving model generalizability in a variety of contexts can be accomplished at a reasonable cost by using weak labels.

#### 4.5 RQ3. Are the Disentangled Representations Truly Invariant?

HATE-WATCH focuses on invariant traits that are essential for hate speech identification across platforms, and achieves excellent causal disentanglement without depending on hate target labels. We do an additional experiment to assess the model’s ability to learn really invariant features. The experiment’s hypothesis is that, as the causal features are common across platforms, there should be significant overlap in the representations of those features if the model can learn invariant features. We trained HATE-WATCH on YouTube and evaluated its generalization on GAB in order to assess this hypothesis. To evaluate representation invariance, we visualized the causal representations (where possible) from 1,000 occurrences per platform, using t-SNE.

HATE-WATCH successfully captures invariant hate speech traits, akin to its fully supervised version, CATCH, according to the t-SNE plots. Among the methods that

highlight the advantages of incorporating causality into hate speech detection frameworks are HATE-WATCH and other causal approaches like CATCH and PEACE, which show a notable advantage in learning invariant representations over language-model based systems like HateBERT. We speculate that the reason HateBERT and other similar models can't learn invariantly or generalize well is because they rely too much on platform-specific details (e.g., HateBERT's training data comes from Reddit, thus the model may have picked up on platform peculiarities while learning from the data).

## 5 Conclusion and Future Work

This paper introduced HATE-WATCH, a novel framework that uses weakly supervised causal disentanglement to effectively navigate the challenging terrain of online hate speech detection. HATE-WATCH effectively disentangles platform dependent features from the platform invariant features of hate, enabling robust detection across various platforms. It sets a new benchmark for flexible and scalable content moderation by successfully reducing the need for large amounts of auxiliary labeled data—that is, the target of hate—through confidence-based reweighting and contrastive regularization. Our results highlight HATE-WATCH's ability to detect hate speech across platforms, leading to the creation of safer online environments.

Future efforts will aim at bridging the gap between different hate expressions and adapting to the evolving landscape of online discourse with reduced manual intervention making HATE-WATCH effectively identify and adapt to new hate speech patterns, resulting in safer communities.

## 6 Ethical Statement

### 6.1 Hate Speech Datasets: Usage and Anonymity

In our research, we have utilized publicly available, well-established datasets, duly citing the relevant sources and adhering to ethical guidelines for their use. We acknowledge the potentially harmful nature of hate speech examples within these datasets, which could be exploited for malicious purposes. Nonetheless, our objective is to enhance the understanding and mitigation of online hate's adverse effects. We have determined that the advantages of employing these real-world examples to explain our research significantly outweigh the associated risks.

### 6.2 Impact Assessment

The development and deployment of hate speech detection systems necessitate comprehensive impact assessments to gauge their societal implications, concerning freedom of expression and the transparency of detection methods.

**Freedom of Expression and Censorship:** Our research is dedicated to creating algorithms capable of identifying and diminishing the presence of harmful language across various platforms, with a keen awareness of the necessity to protect individuals from hate speech while preserving free speech rights. Our methodologies could be applied

to content moderation on social media platforms like X, Facebook, and Reddit to filter out hate speech. Nonetheless, an ethical dilemma arises from the possibility of false positives, where non-hateful content might be mistakenly classified as hate speech, potentially infringing upon legitimate free speech. As such, we caution against the sole reliance on our algorithms for real-world content moderation without the complementary judgment of human annotators to make final decisions.

**Transparency and Fairness in Detection:** Embracing the values of fairness and impartiality, our work is committed to transparently sharing our methods, findings, and inherent limitations, with a continuous goal of enhancing our system. Our dedication to transparency goes beyond merely disclosing our methodologies and results; it encompasses making our decision-making processes clear and comprehensible to ensure ethical practices are followed.

## References

1. Bengio, Y., Courville, A., Vincent, P.: Representation learning: A review and new perspectives. *IEEE Trans. Pattern Analysis Mach. Intell.* (2013)
2. Bourgeade, T., Chiril, P., Benamara, F., Moriceau, V.: What did you learn to hate? a topic-oriented analysis of generalization in hate speech detection. In: *EACL 2023*
3. Buckley, N., Schafer, J.S.: 'censorship-free' platforms: Evaluating content moderation policies and practices of alternative social media (2022)
4. Bühlmann, P.: Invariance, causality and robustness. *Statistical Science* (2020)
5. Caplan, R., Hanson, L., Donovan, J.: Dead reckoning: Navigating content moderation after "fake news" (2018)
6. Carvalho, P., Caled, D., Silva, C., Batista, F., Ribeiro, R.: The expression of hate speech against afro-descendant, roma, and lgbtq+ communities in youtube comments. *Journal of Language Aggression and Conflict* (2023)
7. Caselli, T., Basile, V., Mitrović, J., Granitzer, M.: Hatebert: Retraining bert for abusive language detection in english. *arXiv preprint arXiv:2010.12472* (2020)
8. Cheah, C.S., Wang, C., Ren, H., Zong, X., Cho, H.S., Xue, X.: Covid-19 racism and mental health in chinese american families. *Pediatrics* (2020)
9. Chopra, S., Hadsell, R., LeCun, Y.: Learning a similarity metric discriminatively, with application to face verification. In: *IEEE CVPR'05* (2005)
10. Das, M., Mathew, B., Saha, P., Goyal, P., Mukherjee, A.: Hate speech in online social media. *ACM SIGWEB Newsletter* (2020)
11. Davidson, T., Warmesley, D., Macy, M., Weber, I.: Automated hate speech detection and the problem of offensive language. In: *ICWSM'17* (2017)
12. Döring, N., Mohseni, M.R.: Male dominance and sexism on youtube: results of three content analyses. *Feminist Media Studies* (2019)
13. He, B., Ziems, C., Soni, S., Ramakrishnan, N., Yang, D., Kumar, S.: Racism is a virus: Anti-asian hate and counterspeech in social media during the covid-19 crisis. In: *IEEE/ACM ASONAM* (2021)
14. Jeong, U., Sheth, P., Tahir, A., Alatawi, F., Bernard, H.R., Liu, H.: Exploring platform migration patterns between twitter and mastodon: A user behavior study. *arXiv preprint arXiv:2305.09196* (2023)
15. Jin, Y., Wanner, L., Kadam, V., Shvets, A.: Towards weakly-supervised hate speech classification across datasets. In: *ACL WAOH* (2023)

16. Kikkiseti, D., Mustafa, R.U., Melillo, W., Corizzo, R., Boukouvalas, Z., Gill, J., Japkowicz, N.: Using llms to discover emerging coded antisemitic hate-speech emergence in extremist social media. arXiv preprint arXiv:2401.10841 (2024)
17. Kingma, D.P., Mohamed, S., Jimenez Rezende, D., Welling, M.: Semi-supervised learning with deep generative models. *NeurIPS* (2014)
18. Kullback, S.: *Information theory and statistics*. Courier Corporation (1997)
19. Kumarage, T., Bhattacharjee, A., Garland, J.: Harnessing artificial intelligence to combat online hate: Exploring the challenges and opportunities of large language models in hate speech detection (2024)
20. Kurrek, J., Saleem, H.M., Ruths, D.: Towards a comprehensive taxonomy and large-scale annotated corpus for online slur usage. In: *ACL WAOH* (2020)
21. Lee, R.K.W., Cao, R., Fan, Z., Jiang, J., Chong, W.H.: Disentangling hate in online memes. In: *ACM MM* (2021)
22. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., et al.: Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461 (2019)
23. Liu, Y., Ott, M., Goyal, N., et al.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
24. Locatello, F., Bauer, S., Lucic, M., Rätsch, G., Gelly, S., Schölkopf, B., Bachem, O.: A sober look at the unsupervised learning of disentangled representations and their evaluation. *JMLR* (2020)
25. Macklin, G.: The christchurch attacks: Livestream terror in the viral video age. *CtC Sentinel* (2019)
26. Mansur, Z., Omar, N., Tiun, S.: Twitter hate speech detection: A systematic review of methods, taxonomy analysis, challenges, and opportunities. *IEEE Access* (2023)
27. Mathew, B., Saha, P., Yimam, S.M., Biemann, C., Goyal, P., et al.: Hatexplain: A benchmark dataset for explainable hate speech detection. In: *AAAI* (2021)
28. Meng, Y., Zhang, Y., Huang, J., Xiong, C., Ji, H., Zhang, C., Han, J.: Text classification using label names only: A language model self-training approach. arXiv preprint arXiv:2010.07245 (2020)
29. Pamungkas, E.W., Basile, V., Patti, V.: A joint learning approach with knowledge injection for zero-shot cross-lingual hate speech detection. *Information Processing & Management* (2021)
30. Pereyra, G., Tucker, G., Chorowski, J., Kaiser, Ł., Hinton, G.: Regularizing neural networks by penalizing confident output distributions. arXiv preprint arXiv:1701.06548 (2017)
31. Ramponi, A., Tonelli, S.: Features or spurious artifacts? data-centric baselines for fair and robust hate speech detection. In: *NAACL* (2022)
32. Ramponi, A., Tonelli, S.: Features or spurious artifacts? data-centric baselines for fair and robust hate speech detection. In: *NAACL* (2022)
33. Ren, W., Li, Y., Su, H., Kartchner, D., Mitchell, C., Zhang, C.: Denoising multi-source weak supervision for neural text classification. arXiv preprint arXiv:2010.04582 (2020)
34. Salminen, J., Almerakhi, H., Milenković, M., Jung, S.g., An, J., Kwak, H., Jansen, B.: Anatomy of online hate: developing a taxonomy and machine learning models for identifying and classifying hate in online news media. In: *ICWSM'18* (2018)
35. Sheth, P., Kumarage, T., Moraffah, R., Chadha, A., Liu, H.: Peace: Cross-platform hate speech detection-a causality-guided framework. In: *ECML PKDD* (2023)
36. Sheth, P., Moraffah, R., Kumarage, T.S., Chadha, A., Liu, H.: Causality guided disentanglement for cross-platform hate speech detection. In: *ACM WSDM'24* (2024)
37. Taigman, Y., Yang, M., Ranzato, M., Wolf, L.: Deepface: Closing the gap to human-level performance in face verification. In: *IEEE CVPR* (2014)

38. del Valle-Cano, G., Quijano-Sánchez, L., Liberatore, F., et al.: Socialhaterbert: A dichotomous approach for automatically detecting hate speech on twitter through textual analysis and user profiles. *Expert Systems with Applications* (2023)
39. Wang, Z., Mekala, D., Shang, J.: X-class: Text classification with extremely weak supervision. *arXiv preprint arXiv:2010.12794* (2020)
40. Waseem, Z.: Are you a racist or am i seeing things? annotator influence on hate speech detection on twitter. In: *NLP-CSS* (2016)
41. Yang, C., Zhu, F., Liu, G., Han, J., Hu, S.: Multimodal hate speech detection via cross-domain knowledge transfer. In: *ACM MM* (2022)
42. Yin, W., Agarwal, V., Jiang, A., Zubiaga, A., Sastry, N.: Annobert: Effectively representing multiple annotators' label choices to improve hate speech detection. *arXiv preprint arXiv:2212.10405* (2022)
43. Yu, Y., Zuo, S., Jiang, H., Ren, W., Zhao, T., Zhang, C.: Fine-tuning pre-trained language model with weak supervision: A contrastive-regularized self-training approach. *arXiv preprint arXiv:2010.07835* (2020)
44. Yuan, L., Chen, Y., Cui, G., Gao, H., Zou, F., Cheng, X., Ji, H., Liu, Z., Sun, M.: Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and llms evaluations. *NeurIPS* (2024)