

SHARED CON: Implicit Hate Speech Detection using Shared Semantics

Hyeseon Ahn¹, Youngwook Kim², Jungin Kim¹, and Yo-Sub Han¹

¹Yonsei University, Seoul, Republic of Korea
{hsan, jungin.kim, emmous}@yonsei.ac.kr

²KT, Seoul, Republic of Korea
young-wook.kim@kt.com

Abstract

The ever-growing presence of hate speech on social network services and other online platforms not only fuels online harassment but also presents a growing challenge for hate speech detection. As this task is akin to binary classification, one of the promising approaches for hate speech detection is the utilization of contrastive learning. Recent studies suggest that classifying hateful posts in just a binary manner may not adequately address the nuanced task of detecting implicit hate speech. This challenge is largely due to the subtle nature and context dependency of such pejorative remarks. Previous studies proposed a modified contrastive learning approach equipped with additional aids such as human-written implications or machine-generated augmented data for better implicit hate speech detection. While this approach can potentially enhance the overall performance by its additional data in general, it runs the risk of overfitting as well as heightened cost and time to obtain. These drawbacks serve as motivation for us to design a methodology that is not dependent on human-written or machine-generated augmented data for training. We propose a straightforward, yet effective, clustering-based contrastive learning approach that leverages the shared semantics among the data.¹

1 Introduction

Warning: *this paper contains content that may be offensive or upsetting.*

Detecting hate speech is more challenging compared to other tasks due to its varying perception and definition among individuals, which is reflected in the inconsistencies between datasets (Luo et al., 2023). Implicit hate speech, in particular, often lacks explicit cues such as swear words or insults,

¹Our code is available at <https://github.com/hsannn/sharedcon>.

Shared Semantics:	
White people built everything in this country.	
Posts	Implications
White people built everything in this country.	Whites are superior.
kkk alive and well white brother hood unite ku klux klan.	Minorities are inferior.
if it's white it's right.	Non-white people are wrong.
many whites are waking up to the wonder and pride of being white.	White people should be proud to be white.

Figure 1: An example of shared semantics among posts in the IHC dataset. We often find a representative post among posts sharing similar semantics, which is more accurate than human-written implications. The highlighted portion represents the common meaning shared among these posts.

making its detection even more challenging. Previous studies suggested models for identifying implicit hate speech (Sridhar and Yang, 2022; Lin, 2022). While these approaches exhibit promising performance for the in-dataset scenario, the performance is significantly degraded for the cross-dataset scenario (Wiegand et al., 2019) because of their inconsistent nature. Recent studies have explored methods such as data augmentation (DA) to generate positive pairs (Kim et al., 2022), or creating implicit meanings (*implication*) of hateful posts manually by human (ElSherief et al., 2021). Especially, Kim et al. (2022) showed the importance of the model’s generalization at implicit hate speech detection and proposed a contrastive learning approach that utilizes common implications in the datasets. Although their approach shows good performance for in-dataset, using implications of in-dataset makes it difficult for them to perform well in cross-dataset scenarios. In other words, they

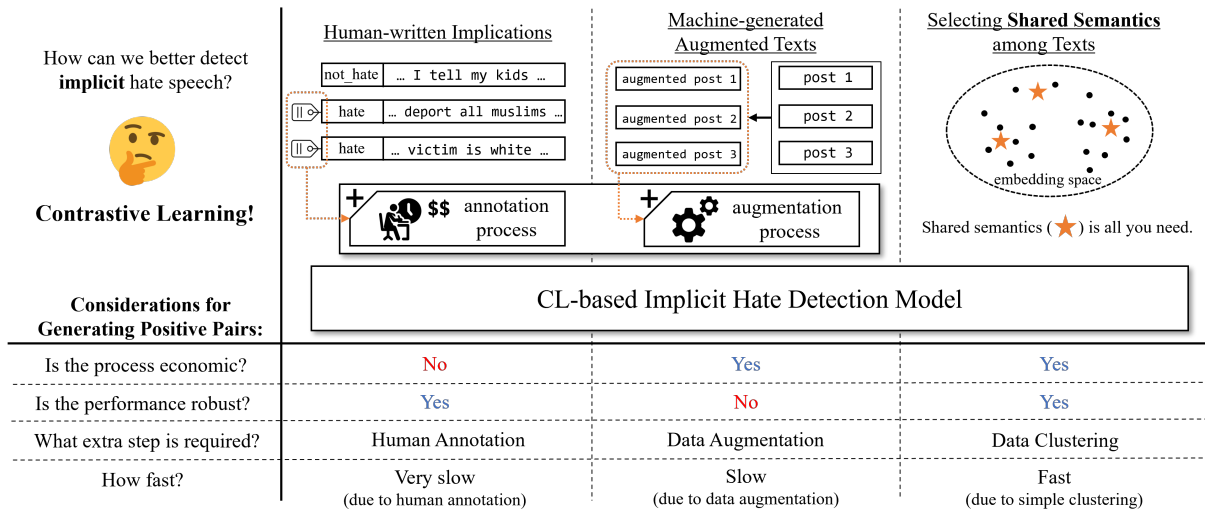


Figure 2: **Middle:** The DA approach does not require special annotations but suffers from lower performance. **Left:** the *implication* approach yields better performances, but its reliance on implication-labeled datasets presents drawbacks in terms of cost and time for creation. **Right:** Our proposed *SharedCon* effectively detects implicit hate speech without relying on human-generated implications, thereby addressing economic concerns and enhancing the robustness of implicit hate speech detection.

are hard to generalize the model in cross-dataset.

In the realm of implicit hate speech detection, the SOTA performance has been achieved by contrastive learning (CL) with human-written implications or machine-generated augmented texts. However, there are some limitations to these approaches: (1) It requires experts to manually generate implications, which is expensive. For example, for an implicit hate speech dataset like IHC (EISherief et al., 2021), the total annotation cost was \$15k and it took a considerable amount of time to complete the dataset construction due to the complexity and scale of the work involved in developing a resource for understanding implicit hate speech. This limitation significantly restricts the scalability of this approach when applied to a wide range of datasets or different domains (Sheth et al., 2023). Furthermore, (2) it is hard to say that DA approach uses the correct positive pairs. The DA strategy merely employs post-augmentation as a positive pair, which lacks semantic consistency across posts. Figure 2 shows the comparison between the existing CL approaches and our methodology.

In order to address implication costs and enhance the existing approaches, we propose *SharedCon*, a simple yet effective technique—leveraging shared semantics from the training data itself without additional data construction steps or extra modules. We extract shared semantics of training datasets by grouping posts and consider the centroid post of each group as the representative semantics shared

by the posts. *SharedCon* is evaluated on three hate speech benchmark datasets in both in-dataset and cross-dataset settings. Experimental results demonstrate that *SharedCon* outperforms our baselines with a 0.43%p average improvement in the in-dataset setting and 1.43%p average improvement in the cross-dataset setting. Furthermore, we compare the performance of *SharedCon* with current SOTA approach (*ImpCon*). The results show that *SharedCon* can deliver comparative (even outperform for some cases) performance, against the current SOTA approach that requires a human-labor-intensive step. Our empirical analysis displays that utilizing shared semantics in the quality-assured sentence embedding space can work as effective as manually produced implications. While the representation of shared semantics works similarly to the implication in Kim et al. (2022), our process provides the flexibility to train the model on various datasets, even those without human-annotated implications. The approach is simple yet effective without extra modules; we cluster the training data and select one post as a core sample of shared semantics, then the selected post plays the role of the anchor in the CL step.

Our main contributions are

- We propose a straightforward and effective technique—leveraging shared semantics from the training data without the need for human-annotated implications.

- We verify the efficacy of shared semantics by enhancing performance across both in-dataset and cross-dataset scenarios.
- We achieve performance that is comparable to and in some cases, even outperforms the current SOTA model, all without the need for human labor-intensive steps.

2 Contrastive Learning with Shared Semantics

While the idea of utilizing shared semantics in the fine-tuning process is promising for implicit hate speech detection, one major drawback is that it requires a dataset with implications, which are often written by humans manually. Therefore, we study an effective way of finding shared semantics without human involvement and improving the cross-dataset performance.

2.1 Shared Semantics

We use the term **shared semantics** to describe the approach for finding shared semantics without human involvement, which we assume that instances of real-world data exhibit inherent shared characteristics. Figure 1 illustrates an example of shared semantics extracted from the IHC (EISherief et al., 2021) dataset; we observe that shared semantics effectively captures the latent meaning of hate posts, which corresponds to implication, a human-generated label representing the implicit meaning of a post.

As depicted in Figure 3, we first extract sentence embeddings of the training data using the existing sentence embedding models. Let us denote a post in the train set and its corresponding label pairs as $\{x_t, y_t\}_{t \in T}$, where T is a total size of a train set. Then, a sentence embedding model constructs high-quality embeddings of x_t . We use the K-means algorithm to cluster the embeddings with the same label. Let C_k be clusters created using K-means and μ_k be a centroid of a cluster C_k , where $k \in [1, 2K]$.² K is a manually adjustable parameter that represents the number of clusters. Each cluster formed in this way is expected to share a common semantic meaning. From this perspective, we consider the centroid posts as the closest sentence from the corresponding cluster centers and define our shared semantics S_k to be

²We use K clusters for each of two labels (hate or not).

$$S_k = \operatorname{argmin}_{x_c \in C_k} \|x_c - \mu_k\|, \quad (1)$$

where x_c is a post in a cluster C_k . We then train the model to bring posts within the same cluster closer together and to push posts from different clusters apart.

2.2 SHARED CON

We propose *SharedCon* that adjusts supervised contrastive learning loss in Eq. (2) by Khosla et al. (2020) for implicit hate speech detection. We denote a pair of input post and its corresponding label as $\{x_i, y_i\}_{i \in I}$, where $I = \{1, \dots, B\}$ is a set of indices in a mini-batch and B is the batch size. As a positive sample x_p for an input post x_i , we train the model to bring x_i and x_p closer together as a positive pair and to push x_i and x_a further apart as a negative pair:

$$\mathcal{L}_{SCL} = \sum_{i=1} \frac{-1}{|\mathcal{P}_i|} \sum_{p \in \mathcal{P}_i} \log \frac{e^{h_i \cdot h_p / \tau}}{\sum_{a \in I \setminus \{i\}} e^{h_i \cdot h_a / \tau}}, \quad (2)$$

where h_i , h_p and h_a are the representations from an encoder for the inputs x_i , x_p and x_a , respectively. \mathcal{P}_i is the set of indices of positive samples and τ is a temperature parameter which is specified in Section 3.3.

A mini-batch is augmented with shared semantics in Eq. (1) per input. Having a mini-batch of B original samples, we set the location of the corresponding shared semantics of the i^{th} ($i \leq B$) original sample x_i as $i + B$. This ensures that each input post in a mini-batch is paired with its positive counterpart, resulting in a total of $2B$ samples for each mini-batch.

We modify Eq. (2) and propose our *SharedCon* loss as follows:

$$\mathcal{L}^{\text{SharedCon}} = \frac{\sum_{p \in \text{SHA}(x_i)} \log \frac{e^{h_i \cdot h_p / \tau}}{\sum_{a \in I \setminus \{i\}} e^{h_i \cdot h_a / \tau}}}{|\text{SHA}(x_i)|}. \quad (3)$$

This supervised contrastive learning method aims to shorten the distance between training samples and their corresponding shared semantics. Step 3 and step 4 in Figure 3 illustrate this process. Here, we consider the scenario where posts from the same cluster appear in a mini-batch. Every post from the same cluster as an input post is

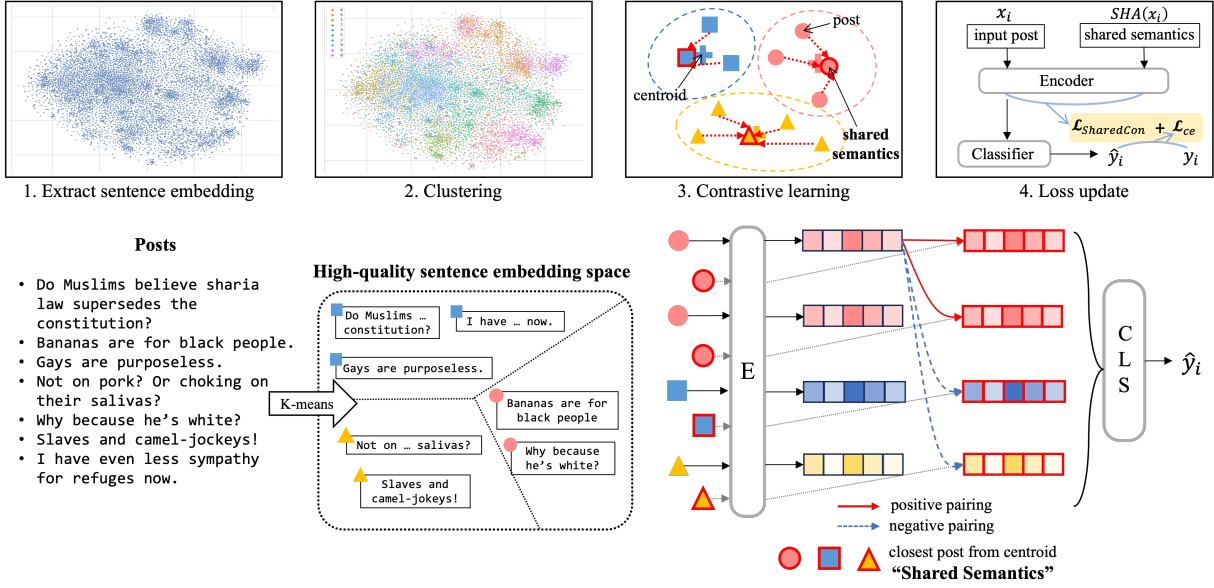


Figure 3: Overall workflow of *SharedCon*. With a centroid sentence (*i.e.*, an anchor) from the clustered sentence embedding space, we train to make sentences within the same cluster closer and sentences from different clusters more distant. ‘E’ denotes an encoder and ‘CLS’ denotes a classifier.

also considered positive. In other words, we allow multiple positives within a mini-batch when members of the same cluster are included.

Suppose that a module $SHA(\cdot)$ constructs our positive set by assigning shared semantics. Then, we arrange the positive sample of x_i as x_p , where $p \in SHA(x_i)$. For $i \leq B$, $SHA(x_i)$ represents the set of shared semantics of the input x_i , but for $i > B$, $SHA(x_i)$ becomes the set of original input posts that come from the same cluster as x_i . Note that we set the location of the original input posts to be in $[1, B]$ and the location of each input’s corresponding positive pair—shared semantics—to be in $[B + 1, 2B]$. We construct our positive set $SHA(x_i)$ as follows:

$$SHA(x_i) = \begin{cases} \{p \mid \mathcal{C}(x_i) = \mathcal{C}(x_p), p \in [B + 1, 2B]\}, & \text{for } i \leq B \\ \{p \mid \mathcal{C}(x_i) = \mathcal{C}(x_p), p \in [1, B]\}, & \text{for } i > B, \end{cases}$$

where $\mathcal{C}(x_i)$ represents the cluster number of x_i . $SHA(x_i)$ serves the role of allowing members within the same cluster to function as positive samples. During training, the model learns the implication of hate speech posts through the *SharedCon* loss in Eq. (3).

The cross-entropy loss is a typical loss function

used in hate speech detection and is defined as follows:

$$\mathcal{L}_{CE} = -\frac{1}{B} \sum_{i=1}^B [y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)], \quad (4)$$

where \hat{y}_i is the model prediction for the x_i and y_i is the ground truth label of x_i .

Since the two losses optimize different aspects, it is necessary to combine them. The contrastive learning loss in Eq. (3) facilitates learning representations that move closer to the centroid as training progresses for better generalization, while the cross-entropy loss in Eq. (4) is used for binary classification itself.

The overall objective is

$$\mathcal{L}_{Overall} = \lambda \mathcal{L}_{CE} + (1 - \lambda) \mathcal{L}_{SharedCon},$$

where λ is a scaling parameter that controls the balance between the cross-entropy loss and the contrastive learning loss.

3 Experimental Results

3.1 Datasets

We conduct multiple experiments utilizing three hate speech datasets reported in Table 1. Note that IHC and SBIC have implications but DynaHate does not.

- **IHC** (ElSherief et al., 2021): a hate speech dataset derived from hate communities on Twitter and their followers, containing information about the target and the implied meaning. The dataset contains 6,346 implicit hate speech out of 22,056 tweets in the United States.
- **SBIC** (Sap et al., 2020): an offensive language dataset that has hierarchical annotations related to social bias, encompassing aspects such as offensiveness, target and implied statements.
- **DynaHate** (Vidgen et al., 2021): a hate speech dataset obtained by a combination of human inputs and model manipulations to deceive the model.

Dataset	Train set	Valid set	Test set
IHC	11,199	3,733	3,734
SBIC	35,504	4,673	4,698
DynaHate	33,006	4,125	4,124

Table 1: The statistical information of three datasets in our experiments.

3.2 Baselines

We consider the following three baselines for comparison.

CE We use the cross-entropy loss (CE) to fine-tune the model, which is a fundamental method for classifying hate speech.

AugCon utilizes *augmented inputs* as positive samples for inputs in the process of contrastive learning. While *AugCon* can be done without any additional need for human annotations (except for inputs and their labels), *AugCon* showed improved performance on some limited settings.

ImpCon utilizes *shared implications* as positive samples for hate speeches in the process of contrastive learning. While *ImpCon* offered consistently superior performance compared to other methods, it necessitates additional human annotation for the implication.

We design a contrastive learning loss that has strengths of both approaches; given that utilizing shared semantics among inputs was effective, we aim to utilize shared semantics for achieving generalization ability (= *ImpCon* strength) without additional human involvements (= *AugCon* strength).

3.3 Implementation Details

For the hyperparameter optimization, we select the learning rate from {5e-6, 1e-5, 2e-5, 3e-5, 5e-5} and the temperature parameter τ from {0.1, 0.3, 0.5}. We also specify λ as a scaling factor from {0.25, 0.5, 0.75} and choose the number K of clusters from {10, 25, 50, 75, 100, 125, 150, 200}. The model archiving the highest validation F1 score is selected as the best-performing one. All experiments are executed using 3 different random seeds and we report the average score of macro-F1 for the results. Unless otherwise stated, we use the SimCSE’s embeddings and the model is trained with the IHC dataset. We perform training on models for a span of 6 epochs employing an NVIDIA RTX 4090 GPU.

3.4 Detection Performance

We verify the effectiveness of our method by evaluating the performance on both in-dataset and cross-dataset scenarios. We train BERT (Devlin et al., 2019) on each of the three datasets and evaluate them on those three respectively. Table 2 demonstrates the outcomes of in-dataset evaluations, with respect to sentence embedding models: SimCSE (Gao et al., 2021). Table 3 presents the evaluation results for cross-dataset scenarios, utilizing SimCSE embeddings as well.

	IHC	SBIC	DYNA
CE	77.7	83.8 [†]	78.8 [†]
+ <i>AugCon</i>	77.4	83.3	77.6
+ <i>ImpCon</i>	78.0 [†]	83.6	-
+ <i>SharedCon</i>	78.5	84.3	79.1

Table 2: In-dataset performance in macro-F1. We report the average of three runs with different random seeds. [†] represents the previous SOTA score, while our score surpassing it is **bolded**.

In the experiments of the in-dataset scenario, we observe a slight improvement of approximately 0.43%p on average over the baselines in Section 3.2. Furthermore, in the cross-dataset setting, particularly for models trained on IHC, there was a significant increase of up to 4.5%p in SBIC evaluation scenario. The noteworthy thing is that the DynaHate dataset does not have implications; our approach selects appropriate representatives of shared semantics from the training dataset and improves the performance.

These results demonstrate the effectiveness of our

Test	Train											
	IHC				SBIC				DYNA			
	CE	Aug	Imp	Ours	CE	Aug	Imp	Ours	CE	Aug	Imp	Ours
IHC	-	-	-	-	59.6	59.7	61.4 [†]	62.4	66.0 [†]	65.6	-	67.0
SBIC	56.8	58.1	60.7 [†]	65.2	-	-	-	-	68.2 [†]	66.3	-	67.9
DYNA	53.1	54.6	57.9 [†]	59.5	60.3	61.2 [†]	61.2 [†]	62.0	-	-	-	-

Table 3: Performance evaluation in macro-F1 for three cross-datasets compared with *AugCon* and *ImpCon*. We report the average of three runs with different random seeds. [†] represents the previous SOTA score, while our score surpassing it is **bolded**.

approach even when compared to the utilization of human-annotated implications. Overall, the experimental results show that *SharedCon* shows competitive or superior performance compared to the existing SOTA without using any human-written implications.

4 Analysis

Our discussion focuses on validating the hypothesis that shared semantics should be identified from the centroid of each cluster, based on a twofold additional scenarios. 1) In Section 4.1, we randomly select a post and allocate it as an anchor; and 2) in Section 4.2, we choose one sentence among the top three or five sentences that are most proximal (in terms of the Euclidean distance) to the centroid of the cluster as an anchor.

4.1 Random Anchors

We examine where shared semantics are located in the embedding space by altering the method of selecting anchor posts. We replace the closest post from the centroid of a cluster in Eq. (1) with $S^* = \text{RANDOM}(x_t)$ such that $\text{RANDOM}(\cdot)$ selects a random post, where $t \in T$ and T is the count of posts within a train set. When a post S^* from a random position is chosen as an anchor for contrastive learning framework, the performance decreases compared to selecting a centroid sentence as an anchor. Figure 4 (Left) demonstrates the superior performance of *SharedCon* compared to the random-anchor-selection approach. The gray graph represents macro-F1 score of *SharedCon* and the violet graph represents the score for random selection. K represents the number of clusters formed in the train set for utilizing the sentence embeddings generated via SimCSE (Gao et al., 2021). As K changes, both selection strategies show minor fluctuations in macro-F1 score, but *SharedCon* consistently achieves good performance. We observed

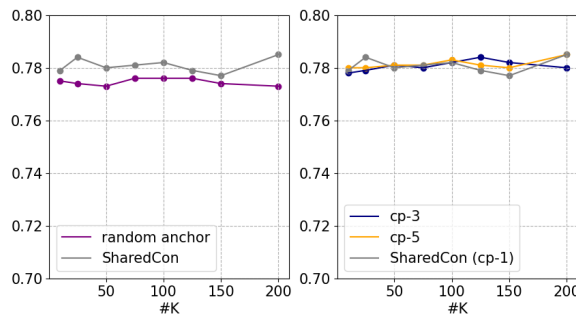


Figure 4: **Left:** The performance when anchors are randomly selected. An overall decrease in performance can be observed in random anchor selection. **Right:** The performance when samples near the centroid are selected as anchors. ‘cp-3 and 5’ refers to the selection of the anchor from the three and five posts most proximal to the centroid, respectively.

that the in-dataset score declined by 1.2%p and the cross-dataset scores decreased by 4.5%p on average in random-anchor experiment setting. These experimental results confirm that shared semantics are located at the center of clusters in a dataset.

4.2 Centroid-proximate Anchors

We conducted additional experiments, applying variations to the previously established positive pairs, to further investigate the methodology for extracting shared semantics. Here, we select the shared semantics in two ways: randomly selecting the shared semantics 1) from the three posts nearest to the centroid of the cluster and 2) from the five posts nearest to the cluster’s centroid. Therefore, we modify Eq. (1) and have

$$S_k = \underset{x_c \in C_k}{\operatorname{argmin}_m} \|x_c - \mu_k\|$$

$$S^* = \text{RANDOM}(S_k),$$

where $\operatorname{argmin}_m \|x_c - \mu_k\|$ denotes the set of top- m posts that minimize the distance of the post x_c

from the cluster centroid μ_k . We set $m \in \{3, 5\}$. Figure 4 (Right) depicts two results of changing the anchor for contrastive learning from the centroid post among the top-3 and top-5 closest posts to the centroid. In the embedding space of SimCSE, we trained the model with the IHC dataset. Both performances of 3-anchors’ (cp-3) and 5-anchors’ (cp-5) are comparable to that of *SharedCon* (cp-1), or occasionally surpass it. This suggests that choosing centroid-around posts might play a role as noise and contribute to the effectiveness of adversarial training. We hypothesize the presence of shared semantics in the center of clusters, but consider the exact location to be an open question for precise positioning. We expect the exploration of alternative methods for selecting representative points beyond the cluster centroids will be reserved for future work. This aspect remains an avenue for further investigation and enhancement in our research. Hence, it can be observed that the centroid sentence performs most effectively in embodying the role of shared semantics.

4.3 Generalizability

While our approach marginally outperforms in-dataset results compared to our baselines, as Kim et al. (2022) tackled, the evaluation of cross-dataset is important to show the generalization ability of a model. Therefore, we conduct a cross-dataset evaluation for the models trained on the three hate speech benchmark datasets: IHC, SBIC, DynaHate (ElSherief et al., 2021; Kim et al., 2023; Vidgen et al., 2021) to check generalizability of our proposed model. Table 3 compares cross-dataset evaluation results of *AugCon*, *ImpCon* and *SharedCon* in SimCSE (Gao et al., 2021) embedding space. We observe that IHC-trained model achieved the highest performance while evaluated in SBIC dataset. We speculate that the reason for the high performance could be attributed to the potential presence of label errors within the dataset. After suspecting potential errors within the IHC dataset, we identified certain noisy implication annotations as shown in Table 4. These labels contribute to confusion in the model, resulting in a decrease in its performance. In other words, it is better to find shared semantics among the given training dataset directly instead of using human-written implications or machine-generated augmented texts.

P1	when you drab the swamp in sacramento!
I1	Immigrants should be deported.
P2	what do you expect? a white kid with air jordan shoes.
I2	Blacks have criminal tendencies
P3	non-white hate crime mob violence comes to finland
I3	Non-white people are violent.

* P is a post and I is its implication

Table 4: Examples of wrong implications written by human-annotators in the IHC dataset. Such errors can often occur in the dataset.

4.4 Representation by Sentence Embedding Models

We assume that the space of the sentence embedding model already effectively reflects semantic similarity. In order to validate our assumption, we also implemented an extra experiment by changing the sentence embedding model. Figure 5 shows that *SharedCon* consistently demonstrates strong performance regardless of variations in the sentence embedding model.

Sentence embedding model We leverage three different sentence embedding models: SBERT (Reimers and Gurevych, 2019), SimCSE (Gao et al., 2021) and Angle (Li and Li, 2023)

- **SBERT** is a modified BERT using siamese / triplet structures, excels in semantic similarity and transfer learning.
- **SimCSE** enhanced sentence embeddings through contrastive learning, improving semantic understanding in textual tasks via unsupervised and supervised methods.
- **Angle** addressed gradient issues in cosine-based methods across diverse text types and outperformed prior models in semantic similarity tasks.

By assuming that the space represented by a sentence embedding model embodies a comprehensive semantic structure, we discerned shared semantics within this framework. We aim to empirically substantiate the validity of high-quality representation space. Figure 5 shows the results for SBERT, SimCSE and Angle consistently exhibited favorable performance without significant deviations. In this section, we train the model with the IHC dataset and present the average results from five random

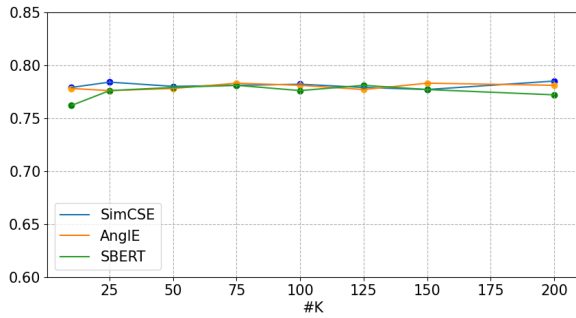


Figure 5: We compare F1 score of three different sentence embedding models. While the range of K is from 10 to 200, the performance using the three models is similar, but a slight advantage can be observed in SimCSE and AnglE compared to SBERT.

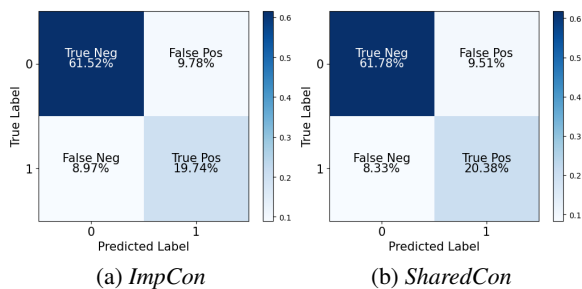


Figure 6: Confusion matrices for the model predictions on IHC test set. We evaluate the predictions made by two BERT models trained on IHC train set with (a) *ImpCon* and (b) *SharedCon*.

seeds. The consistency in the performance of these three models demonstrates the embedding invariance of our approach. This indicates that *SharedCon* can achieve fair performance regardless of which model is used.

The effectiveness of our methodology lies in extracting implications within a high-quality representation (embedding) space. In other words, the sentence embedding model provides a reasonably good representation space for finding shared semantics and our approach of identifying positive pairs on this basis makes sense.

4.5 Error Analysis

We analyze the incorrect predictions of our approach and *ImpCon* for the IHC dataset. *SharedCon* had 355 false positives and 311 false negatives while *ImpCon* had 365 and 335 respectively. As depicted in the confusion matrices in Figure 6, *SharedCon* has 0.64%p less false negatives and 0.27%p less false positives compared to *ImpCon*. Table 5 shows some examples of errors both *ImpCon* and *SharedCon* made. For *SharedCon*, out

of a total of 666 incorrections in both false positive and false negative, we observed that there are 22 broken sentence posts for both in false positive and false negative. That is, both methodologies encounter difficulties with the broken sentences, while *ShareCon* is comparatively robust than *ImpCon* with grammatically correct sentences. *ImpCon* is more vulnerable to broken sentences than *SharedCon*; clearly, broken sentences impact the model’s prediction errors. Since broken sentences and obfuscated posts are commonly observed in recent social media platforms, effectively detecting these aspects is an important avenue for future research. We provide samples of broken sentences that both methodologies struggle with at Table 5.

Regarding the misclassification by *SharedCon*, we point out errors within the dataset itself. We often observe errors in human-written implications in the dataset; Table 4 shows some examples that capture subtly incorrect implications. These examples illustrate instances where responses are skewed towards negative connotations or mislabeled altogether, irrespective of context. For a very basic analysis to identify the simple label error ratio, we randomly sampled 100 posts from the train set of IHC and found that approximately 27% of them were misclassified. Out of these misclassifications, about 7 posts were ambiguous and the remaining 20 posts were incorrectly classified. We assume that these incorrect implications can cause further performance degradation.

5 Related Work

Hate Speech Detection There are several datasets on explicit hate speech with different purposes (Warner and Hirschberg, 2012; Davidson et al., 2017; Founta et al., 2018; Basile et al., 2019). Davidson et al. (2017) addressed distinguishing hate speech from offensive language by using a crowd-sourced hate speech lexicon. Implicit hate speech, on the other hand, characterized by its subtle linguistic expressions, presents a more complex obstacle for detection compared to its explicit hate speech (Ocampo et al., 2023). Recently, researchers have started considering implicit hatred or abusiveness. Han and Tsvetkov (2020) introduced a novel generic framework, MACAS, for online abusive language detection across different domains. Wiegand et al. (2021) addressed the challenges in detecting implicitly abusive language, highlighting the limitations of existing datasets and

Label: 1	Post1	Post2	Post3
<i>ImpCon</i> : 0 <i>SharedCon</i> : 1	i happen to prefer the white race not partially covered in non-white ink.	white power and proud !!!	it is muslim law to kill gay pple ?
Label: 1	Post1	Post2	Post3
<i>ImpCon</i> : 0 <i>SharedCon</i> : 0	to say " white lives matter " " is pure " " racist hate " " . anti-racist is a code word for #antiwhite . " "	rt " : super cool german officer puts " " neo-nazi " " stickers in police van more" ; "	latest track " ethnostate 2088 " " "

Table 5: Error examples from IHC train dataset. We compare prediction of *ImpCon* and *SharedCon*. Label 1 denotes that ground truth label is ‘implicit hate’.

proposed a divide-and-conquer strategy to categorize subtypes of implicit abuse. In order to identify implicit hate speech [ElSherief et al. \(2021\)](#) provided a benchmark dataset to understand implicit hate speech beyond explicit forms. Leveraging the implicit hate speech dataset, [Kim et al. \(2022\)](#) introduced *ImpCon*, a contrastive learning method that addressed the challenge of model generalization in implicit hate speech detection. In recent years, there have been some research proposing model training strategy ([Jafari et al., 2023](#)) for implicit hate speech detection, as well as suggesting new architectures for models ([Ghosh et al., 2023](#)).

Contrastive learning Recently, contrastive learning has become a prominent tool in unsupervised representation learning in various domains. SimCLR ([Chen et al., 2020](#)) randomly selects two augmented examples from a mini-batch as a positive pair, while the other augmented examples within the mini-batch are treated as negative pairs. [Khosla et al. \(2020\)](#) proposed to select positive pairs and negative pairs by considering anchor points. In the natural language processing field, SimCSE ([Gao et al., 2021](#)) selected positive pairs by using annotated pairs from NLI datasets, where entailment pairs are considered as positives and contradiction pairs used as hard negatives. [Gunel et al. \(2021\)](#) suggested using contrastive learning as a fine-tuning tool to capture similarities between examples from the same classes and contrast them with examples from other classes to improve generalization. [Suresh and Ong \(2021\)](#) performed well on fine-grained text classification tasks by adaptively weighting the relationships between classes to distinguish between more difficult samples. Recently, [Lu et al. \(2023\)](#) proposed dual contrastive learning approach for detecting hate speech. We also used contrastive learning to fine-tune model performance.

Clustering for contrastive learning [Yang et al. \(2023\)](#) performs cluster-level contrastive learning to incorporate measurable emotion prototypes. [Gao et al. \(2022\)](#) proposed a weakly supervised contrastive learning method that allows us to consider multiple positives and multiple negatives. These are prototype-based clustering methods that avoid semantically related events being pulled apart. We suggest selecting multiple positives by clustering to use contrastive learning in the fine-tuning process for shared semantics without human-written implications.

6 Conclusions

We have enhanced the existing implicit hate speech detection approach, which previously relied on human-generated implications. This improvement involves clustering sentence embeddings and utilizing the centroid sentence of each cluster as positive examples during the contrastive learning process. The experimental results have demonstrated that posts in the dataset have shared semantics and our approach of using them as positive samples for contrastive learning works well. Particularly, our methodology has shown competitive or superior performance improvement of fine-tuned models in cross-dataset experiments. To sum up, our approach has enhanced model generalization without human-written implications, providing a better solution for implicit hate speech detection in terms of cost.

Limitations

While our approach reduces possible errors of wrong implications written by humans, it might be possible that a centroid sentence itself might be labeled wrong in the first place. However, it is less likely for a centroid sentence to have a wrong label since it has similar semantics compared with the

other sentences in the same cluster. Another issue would be a selection process of shared-semantic sentences. Certainly, there could be a better way, but our approach—selecting a centroid sentence in a cluster—constantly shows good and stable performance compared with several other heuristics.

Ethical Consideration

Generalization Ability Enhancement Our work focuses on enhancing the generalization ability to train a hate speech detector that can consistently adapt to new forms of hate speech. Specifically, by training a model with improved performance, we aim to equip it with the capability to respond to emerging manifestations of hate speech even on the same dataset. This emphasis on enhanced generalization entails training a model that demonstrates superior performance when presented with the same dataset, enabling it to effectively address novel expressions of hate speech.

Annotator Exposure Mitigation In the past, state-of-the-art (SOTA) models exposed annotators to the risk of hate speech exposure during the process of adding implications as comments to secure such capabilities. However, our methodology has achieved comparable performance without such risks.

Risks and Potential Misuse The enhancement of these capabilities represents a positive development in the field of hate speech detection; however, it also inherently carries risks. If the improved generalization ability is exploited, it could be utilized to generate new forms of hate speech or conceal hateful behaviors. It is crucial to be vigilant about these potential risks and to carefully evaluate the use of the model and its outcomes.

Acknowledgements

This research was supported by the NRF grant (RS-2023-00208094) and the AI Graduate School Program at Yonsei University (No. RS-2020-II201361) funded by the Korean government (MSIT). Han is a corresponding author.

References

Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. SemEval-2019 task 5: Multilingual detection

of hate speech against immigrants and women in Twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning*, pages 1597–1607.

Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. *Proceedings of the International AAAI Conference on Web and Social Media*, 11(1):512–515.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.

Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. 2021. Latent hatred: A benchmark for understanding implicit hate speech. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*. Association for Computational Linguistics.

Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of Twitter abusive behavior. In *Proceedings of the 12th International AAAI Conference on Web and Social Media*, pages 491–500.

Jun Gao, Wei Wang, Changlong Yu, Huan Zhao, Wilfred Ng, and Ruifeng Xu. 2022. Improving event representation via simultaneous weakly supervised contrastive learning and clustering. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 3036–3049. Association for Computational Linguistics.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6894–6910. Association for Computational Linguistics.

Sreyan Ghosh, Manan Suri, Purva Chiniya, Utkarsh Tyagi, Sonal Kumar, and Dinesh Manocha. 2023.

- Cosyn: Detecting implicit hate speech in online conversations using a context synergized hyperbolic network. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 6159–6173.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Veselin Stoyanov. 2021. Supervised contrastive learning for pre-trained language model fine-tuning. In *International Conference on Learning Representations*.
- Xiaochuang Han and Yulia Tsvetkov. 2020. Fortifying toxic speech detectors against veiled toxicity. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7732–7739.
- Amir Reza Jafari, Guanlin Li, Praboda Rajapaksha, Reza Farahbakhsh, and Noël Crespi. 2023. Fine-grained emotions influence on implicit hate speech detection. *IEEE Access*, 11:105330–105343.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. In *Advances in Neural Information Processing Systems*, pages 18661–18673.
- Youngwook Kim, Shinwoo Park, and Yo-Sub Han. 2022. Generalizable implicit hate speech detection using contrastive learning. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6667–6679.
- Youngwook Kim, Shinwoo Park, Youngsoo Namgoong, and Yo-Sub Han. 2023. Conprompt: Pre-training a language model with machine-generated data for implicit hate speech detection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10964–10980.
- Xianming Li and Jing Li. 2023. Angle-optimized text embeddings. *arXiv preprint arXiv:2309.12871*.
- Jessica Lin. 2022. Leveraging world knowledge in implicit hate speech detection. *arXiv preprint arXiv:2212.14100*.
- Junyu Lu, Hongfei Lin, Xiaokun Zhang, Zhaoqing Li, Tongyue Zhang, Linlin Zong, Fenglong Ma, and Bo Xu. 2023. Hate speech detection via dual contrastive learning. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:2787–2795.
- Chu Luo, Rohan Bhambhoria, Samuel Dahan, and Xiaodan Zhu. 2023. [Legally enforceable hate speech detection for public forums](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Singapore. Association for Computational Linguistics.
- Nicolas Ocampo, Elena Cabrio, and Serena Villata. 2023. Playing the part of the sharp bully: Generating adversarial examples for implicit hate speech detection. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 2758–2772. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490.
- Paaras Sheth, Tharindu Kumarage, Raha Moraffah, Aman Chadha, and Huan Liu. 2023. Peace: Cross-platform hate speech detection - a causality-guided framework. In *Machine Learning and Knowledge Discovery in Databases: Research Track*, pages 559–575.
- Rohit Sridhar and Diyi Yang. 2022. [Explaining toxic text via knowledge enhanced text generation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 811–826, Seattle, United States. Association for Computational Linguistics.
- Varsha Suresh and Desmond Ong. 2021. Not all negatives are equal: Label-aware contrastive loss for fine-grained text classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4381–4394.
- Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021. Learning from the worst: Dynamically generated datasets to improve online hate detection. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1667–1682.
- William Warner and Julia Hirschberg. 2012. Detecting hate speech on the world wide web. In *Proceedings of the Second Workshop on Language in Social Media*, page 19–26.
- Michael Wiegand, Josef Ruppenhofer, and Elisabeth Eder. 2021. Implicitly abusive language – what does it actually look like and why are we not getting there? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 576–587.
- Michael Wiegand, Josef Ruppenhofer, and Thomas Kleinbauer. 2019. Detection of Abusive Language: the Problem of Biased Datasets. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 602–608.

Kailai Yang, Tianlin Zhang, Hassan Alhuzali, and Sophia Ananiadou. 2023. Cluster-level contrastive learning for emotion recognition in conversations. *IEEE Trans. Affect. Comput.*, 14(4):3269–3280.