

# Understanding and Detecting Hateful Content using Contrastive Learning

Felipe González-Pizarro<sup>♦\*</sup>, and Savvas Zannettou<sup>♦▲</sup>

<sup>♦</sup>University of British Columbia, <sup>♦</sup>TU Delft, <sup>▲</sup>Max Planck Institute for Informatics

felipegp@cs.ubc.ca, s.zannettou@tudelft.nl

## Abstract

The spread of hate speech and hateful imagery on the Web is a significant problem that needs to be mitigated to improve our Web experience. This work contributes to research efforts to detect and understand hateful content on the Web by undertaking a multimodal analysis of Antisemitism and Islamophobia on 4chan’s /pol/ using OpenAI’s CLIP. This large pre-trained model uses the Contrastive Learning paradigm. We devise a methodology to identify a set of Antisemitic and Islamophobic hateful textual phrases using Google’s Perspective API and manual annotations. Then, we use OpenAI’s CLIP to identify images that are highly similar to our Antisemitic/Islamophobic textual phrases. By running our methodology on a dataset that includes 66M posts and 5.8M images shared on 4chan’s /pol/ for 18 months, we detect 573,513 posts containing 92K Antisemitic/Islamophobic images and 246K posts that include 420 hateful phrases. Among other things, we find that we can use OpenAI’s CLIP model to detect hateful content with an accuracy score of 0.84 (F1 score = 0.58). Also, we find that Antisemitic/Islamophobic imagery is shared in 2x more posts on 4chan’s /pol/ compared to Antisemitic/Islamophobic textual phrases, highlighting the need to design more tools for detecting hateful imagery. Finally, we make publicly available a dataset of 420 Antisemitic/Islamophobic phrases and 92K images that can assist researchers in further understanding Antisemitism/Islamophobia and developing more accurate hate speech detection models.

## 1 Introduction

The spread of hateful content on the Web is an everlasting and vital issue that adversely affects society. The problem of hateful content is longstanding on the Web for various reasons. First, there is no scientific consensus on what constitutes hateful content (i.e., no definition of what hate speech is) [56]. Second, the problem is complex since hateful content can spread across various modalities (e.g., text, images, videos, etc.), and we still lack automated techniques to detect hateful content with acceptable and generalizable performance [4]. Third, we lack moderation tools to proactively prevent the spread of hateful content on the Web [35]. This work focuses on assisting the community in addressing the issue of the lack of tools to detect hateful content across multiple modalities.

Most of the research efforts in the space of detecting hateful content focus on designing and training machine learning models that are specifically tailored towards detecting specific instances of hateful content (e.g., hate speech on text or particular cases of hateful imagery). Some examples of such efforts include Google’s Perspective API [46] and the HateSonar classifier [17] that aim to detect toxic and offensive text. Other methods aim to detect instances of hateful imagery like Antisemitic images [71] or hateful memes [33, 69]. These efforts and tools are essential and valuable, however, they rely on human-annotated datasets that are expensive to create, and therefore they are also small. At the same time, these datasets focus on specific modalities (i.e., text or images in isolation). All these drawbacks limit their broad applicability.

The lack of large-scale annotated datasets for solving problems like hate speech motivated the research community to start developing techniques that do not rely on annotated datasets (a paradigm known as *unsupervised learning*). Over the past years, the research community released large-scale models that depend on huge unlabeled datasets such as OpenAI’s GPT-3 [6], Google’s BERT [18], OpenAI’s CLIP [48], etc. These models are trained on large-scale datasets and usually can capture general knowledge extracted from the datasets that can be valuable for performing classification tasks that the model was not explicitly trained on.

Motivated by these recent advancements on large-scale pre-trained machine learning models, in this work, we investigate how we can use such models to detect hateful content on the Web across multiple modalities (i.e., text and images). Specifically, we focus on OpenAI’s CLIP model because it helps us capture content similarity across modalities (i.e., measure similarity between text and images). To achieve this, CLIP leverages a paradigm known as Contrastive Learning; the main idea is that the model maps text and images to a high-dimensional vector space and is trained in such a way that similar text/images are mapped closer to this vector space (for more details see Section 2).

**Focus & Research Questions.** This work focuses on understanding the spread of Antisemitic/Islamophobic content on 4chan’s /pol/ board. We concentrate on hateful content targeted towards these two demographics mainly because previous work indicates that 4chan’s /pol/ is known for disseminating Antisemitic/Islamophobic content [47, 71]. Specifically, we focus on shedding light on the following research questions:

\*Work done during an internship at Max Planck Institute for Informatics.

- **RQ1:** Can large pre-trained models that leverage the Contrastive Learning paradigm, like OpenAI’s CLIP, identify hateful content with acceptable performance?
- **RQ2:** How prevalent is Antisemitic/Islamophobic imagery and textual hate speech on 4chan’s /pol/?

To answer these research questions, we obtain all the posts and images shared on 4chan’s /pol/ between July 1, 2016, and December 31, 2017, ultimately collecting 66M textual posts and 5.8M images. Then, we leverage the Perspective API and manual annotations to construct a dataset of 420 Antisemitic and Islamophobic textual phrases. We retrieve 246K posts that include any of our 420 hateful phrases. Finally, we use OpenAI’s CLIP to detect Antisemitic/Islamophobic images when provided as input the above-mentioned hateful phrases and all images shared on 4chan’s /pol/; we find 92K images.

**Contributions & Main Findings.** Our work makes the following contributions/main findings:

- We investigate whether large pre-trained models based on Contrastive Learning can assist in detecting hateful imagery. We find that large pre-trained models like OpenAI’s CLIP [48] can detect Antisemitic/Islamophobic imagery with 0.84, 0.47, 0.80, 0.58, accuracy, precision, recall, and F1 score, respectively (**RQ1**).
- We find that on 4chan’s /pol/ Antisemitic/Islamophobic imagery appears in 2x more posts compared to Antisemitic/Islamophobic textual hateful content. This finding highlights the need for the development and use of multimodal hate speech detectors for understanding and mitigating the problem (**RQ2**).
- We make publicly available (upon acceptance of the manuscript) a large dataset of Antisemitic/Islamophobic posts, phrases, and images shared on 4chan’s /pol/. The released dataset can assist researchers in future work focusing on detecting and understanding the spread of hateful content on the Web across multiple modalities (i.e., text and images).

**Ethical Considerations.** We emphasize that we rely entirely on publicly available and anonymous data shared on 4chan’s /pol/. The authors of this work did all manual annotations performed in this study, hence there are no other individuals exposed to hateful content. Also, we follow standard ethical guidelines [51] like reporting our results on aggregate and not attempting to deanonymize users.

**Disclaimer.** This manuscript contains Antisemitic and Islamophobic textual and graphic elements that are offensive and are likely to disturb the reader.

## 2 Background

This section provides background information on Contrastive Learning and OpenAI’s CLIP model, as well as details on Google’s Perspective API.

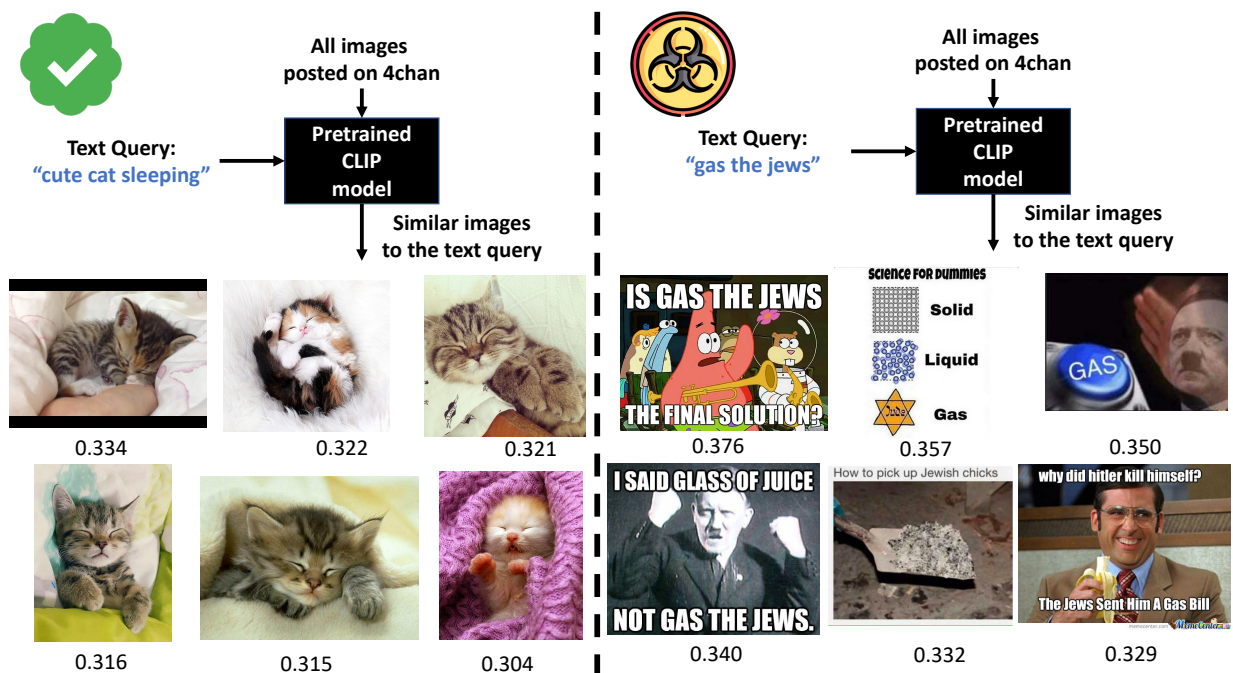
**Contrastive Learning.** To understand Contrastive Learning, it is essential to grasp its differences compared to traditional Machine/Deep Learning (ML/DL) classifiers. Traditional ML/DL classifiers take as an input a set of unlabeled data, each accompanied with a class, and aim to predict the class from the unlabeled data, a paradigm known as supervised learning [7]. On the other hand, Contrastive Learning is an unsupervised technique, meaning that there is no need to have labeled data (i.e., classes). The main idea behind Contrastive Learning is that you train a model that relies entirely on unlabeled data, and the model learns general features from the dataset by teaching it which input samples are similar/different to each other [28]. In other words, Contrastive Learning relies on a set of unlabeled data samples with additional information on which of these samples are similar to each other. Contrastive Learning is becoming increasingly popular in the research community with several applications on visual representations [12, 34], textual representations [22, 23, 65], graph representations [30, 66], and multimodal (i.e., text/images, images/videos, etc.) representations [19, 48, 67, 72].

**OpenAI’s CLIP.** OpenAI recently released a model called Contrastive Language-Image Pre-training (CLIP) [48] that leverages Contrastive Learning to generate representations across text and images. The model relies on a text encoder and an image encoder that maps text and images to a high-dimensional vector space. Subsequently, the model is trained to minimize the cosine distance between similar text/image pairs. To train CLIP, OpenAI created a huge dataset that consists of 400M pairs of text/images collected from various Web sources and covers an extensive set of visual concepts.<sup>1</sup> By training CLIP with this vast dataset, the model learns general visual representations and how these representations are described using natural language, which results in the model obtaining general knowledge in various topics (e.g., identify persons, objects, etc.).

In this work, we leverage the CLIP model to extract representations for our 4chan textual/image datasets and assess the similarity between the text and image representations. The main idea is that by providing to the pre-trained CLIP model a set of hateful text phrases, we will be able to identify a set of hateful images that are highly similar to the hateful text query. To demonstrate CLIP’s potential in discovering hateful imagery from hateful text-based queries, we show an example from our 4chan dataset in Fig. 1. On the left side, we show examples of images that are highly similar to a benign text query like “cute cat sleeping,” while on the right, we show examples of images that are similar to an antisemitic and toxic phrase (“gas the jews”)<sup>2</sup>. For each image, we report the cosine similarity between the representation obtained from the text query and the representation of the image in our dataset. This example shows that the CLIP model can detect objects in images (i.e., cats) and provide relevant images to the queries (i.e., the cats are indeed sleeping according to the query). Further-

<sup>1</sup>The exact methodology for creating this dataset was not made publicly available by OpenAI.

<sup>2</sup>In this work, we treat an image as similar to the text phrase if it has a cosine similarity of 0.3 or higher (see Section 4.2).



**Figure 1:** Example of text queries and images that are similar to the text queries on 4chan (i.e., cosine similarity between the text CLIP-representation and the image CLIP-representation equals to 0.3 or more). On the left side, we show a benign text query (“cute cat sleeping”), while on the right we show the results for a toxic and antisemitic query (“gas the jews”).

more, by looking at the images for the toxic query, we observe that CLIP can identify harmful images based on the query and can link historical persons to it (e.g., the textual input does not mention Adolf Hitler; however, the model knows that Hitler was responsible for the holocaust). Also, CLIP can detect images that share hateful ideology by adding text on memes (i.e., CLIP also performs Optical Character Recognition and can correlate that text with the text-based query). Overall, this example shows the predictive power of the CLIP model in detecting hateful imagery from hateful text phrases.

**Google’s Perspective API.** As a first step towards identifying hateful phrases, we use Google’s Perspective API [46], which provides a set of Machine Learning models for identifying how rude/aggressive/hateful a comment is. We use the Perspective API for identifying hateful text mainly because it outperforms other publicly available hate speech classifiers like HateSonar [17, 70]. This work focuses on the SEVERE\_TOXICITY model available from Perspective API because it is more robust to positive uses of curse words, and it is a production-ready model. The SEVERE\_TOXICITY model returns a score between 0 and 1, which can be interpreted as the probability of the text being rude and toxic.

### 3 Dataset

Our work focuses on 4chan, particularly the Politically Incorrect board (/pol/). /pol/ is the main board for discussing world events and politics and is infamous for the spread of conspiracy theories [58, 68] and racist/hateful content [31, 71]. We collect the data about posts on /pol/ using the publicly available dataset released by Papasavva et al. [45]; the dataset in-

cludes textual data about 134.5M posts shared on /pol/ between June 2016 and November 2019. Our work focuses on the period between July 1, 2016, and December 31, 2017 (to match the time period of the image dataset mentioned below), including 66,383,955 posts. We complement the above dataset with the image dataset collected by Zannettou et al. [71]. The dataset includes 5,859,439 images shared alongside /pol/ posts between July 1, 2016, and December 31, 2017.

Overall, our dataset comprises all textual and image activity on /pol/ between July 1, 2016, and December 31, 2017, including 66M posts and 5.8M images.

## 4 Methodology

This section describes our methodology for detecting hateful text phrases and hateful imagery, focusing on Antisemitic and Islamophobic content.

### 4.1 Identifying Antisemitic and Islamophobic phrases

Here, our goal is to identify a set of phrases that are Antisemitic/Islamophobic. To do this, we follow a multi-step semi-automated methodology. First, we use the SEVERE\_TOXICITY scores from the Perspective API to identify posts that are toxic/offensive without considering the target (e.g., if it is antisemitic). Specifically, we consider all posts that have a score of 0.8 or more as toxic, following the methodology by Ribeiro et al. [50]. Out of the 66M posts in our dataset, we find 4.5M (6.7%) toxic posts.

Having extracted a set of toxic posts from 4chan’s /pol/, we then aim to identify the main targets of hate speech on /pol/

by extracting the top keywords. To do this, we preprocess the data to remove HTML tags, stop words, and URLs, and then we create a term frequency-inverse document frequency array (TF-IDF). Next, we manually inspect the top 200 words based on their TF-IDF values and identify the words related to Jews or Muslims. As a result, we find seven keywords: “jews,” “kike,” “jew,” “kikes,” “jewish,” “muslims,” and “muslim.” Then, based on these keywords, we filter the toxic posts obtained from the previous step, hence getting a set of 336K posts with a SEVERE\_TOXICITY score of 0.8 and include at least one of the seven keywords.

Since our goal is to create a set of Antisemitic/Islamophobic phrases, we need to break down the toxic 4chan posts into sentences and then identify the ones that are Antisemitic/Islamophobic. To do this, we apply a sentence tokenizer [41] on the 336K posts, obtaining 976K sentences. To identify common phrases used on 4chan’s /pol/, we apply WordNet lemmatization [42], excluding all sentences that appear less than five times. We obtain 4,582 unique common phrases; not all of these sentences are Antisemitic/Islamophobic.

Identifying whether a phrase is Antisemitic/Islamophobic is not a straightforward task and can not be easily automated. Therefore, we use manual annotation on the 4,582 common phrases to annotate the common phrases as Antisemitic/Islamophobic or irrelevant. Two authors of this paper independently annotated the 4.5K common phrases. We discard long phrases (over seven words) during the annotation since our preliminary experiments showed that OpenAI’s CLIP returns a considerable amount of false positives when provided with long text queries. We also consider as irrelevant phrases that target multiple demographic groups (e.g., hateful towards Muslims and Jews like “fuck jews and muslims” or hateful towards African Americans and Jews like “fuck niggers and jews”). The two annotators agreed on 91% of the annotations with a Cohen’s Kappa score of 0.69, which indicates a substantial agreement [36]. After the independent annotations, the two annotators discussed the disagreements to come up with a final annotation on whether a phrase is Antisemitic/Islamophobic or irrelevant. After our annotation, we find 326 Antisemitic and 94 Islamophobic phrases. The list of the Antisemitic/Islamophobic phrases is publicly available [3].

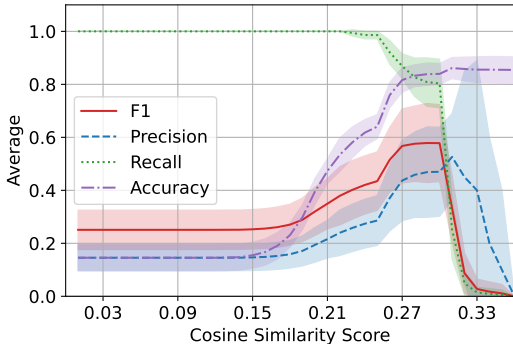
Finally, we search for these Antisemitic/Islamophobic phrases on the entire dataset. We extract all posts that include any of the Antisemitic/Islamophobic phrases (*Textual* dataset), finding 247K posts. Note that we remove 864 posts that contain both Antisemitic and Islamophobic phrases. Overall, we find 209K Antisemitic posts and 37K Islamophobic posts (see Table 1).

## 4.2 Identifying Antisemitic and Islamophobic images

Our goal is to identify Antisemitic and Islamophobic imagery using the pre-trained CLIP model [48]. To do this, we encode all images in our dataset using the image encoder on the CLIP model, hence obtaining a high-dimensional vector for each image. Also, we encode all the Antisemitic/Islamophobic

Dataset	Textual		Visual	
	# Phrases	# Posts	# Images	# Posts
<b>Antisemitism</b>	326	209,224	69,610	472,048
<b>Islamophobia</b>	94	37,354	22,519	101,465
<b>Total</b>	420	246,578	92,129	573,513

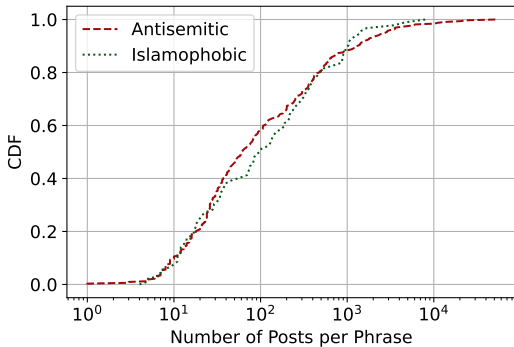
**Table 1:** Overview of our Antisemitism/Islamophobia Textual and Visual datasets. The number of phrases is based on lemmatized versions, and the number of images is based on the unique pHash values.



**Figure 2:** Performance of the CLIP model in identifying Antisemitic/Islamophobic imagery for varying cosine similarity thresholds. The lines refer to the average metric for ten random phrases (2K images), while the area refers to the standard deviation across the ten phrases.

phrases (extracted from the previous step), using the text encoder on the CLIP model, obtaining a vector for each phrase. Then, we calculate all the cosine similarities between the image and text vectors, which allows us to assess the similarity between the phrases and the images. The main idea is that by comparing a hateful phrase to all the images, images with a high cosine similarity score will also be hateful. To identify a suitable cosine similarity threshold where we treat a text and an image similarly, we perform a manual annotation process.

First, we extract a random sample of ten Antisemitic/Islamophobic phrases (eight Antisemitic and two Islamophobic to match the percentage of Antisemitic/Islamophobic phrases in our dataset). Then, we extract a random sample of 200 images for each phrase while ensuring that the images cover the whole spectrum of cosine similarity scores. Specifically, we extract 50 random images with cosine similarity scores for each of the following ranges: [0.0, 0.20), [0.2, 0.25), [0.25, 0.3), [0.3, 0.4]. To select these ranges, we plot the CDF of all cosine similarity scores obtained by comparing the ten randomly selected phrases and all the images in our dataset (we omit the figure due to space constraints). We find that 40% of the scores are below 0.2, and we expect these images to be entirely irrelevant for the phrase. To verify this, we select the [0.0, 0.20) range. Additionally, we select the [0.2, 0.25) because it has a considerable percentage of the scores (50%), and we expect that the images will not be very similar again. Finally, we select the [0.25, 0.3) and [0.3-0.4] ranges because we expect that the ideal threshold is somewhere in these two ranges, and devoting half of the



**Figure 3:** CDF of the number of Antisemitic/Islamophobic posts containing each phrase.

selected images in these ranges will help us identify a suitable threshold.

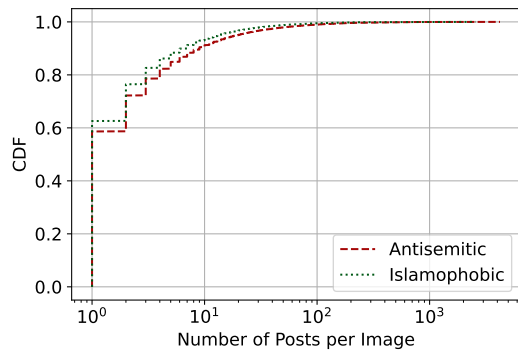
Then, two authors of this paper independently annotated the 2K images to identify which are Antisemitic/Islamophobic or irrelevant. The annotators agreed on 94% of the annotations with a Cohen’s Kappa score of 0.75, which indicates a substantial agreement. Again, the two annotators solved the disagreements by discussing the images and deciding a final annotation on whether the image is Antisemitic/Islamophobic or irrelevant.

Having constructed a ground truth dataset of Antisemitic and Islamophobic imagery, we then find the best performing cosine similarity threshold. We vary the cosine similarity threshold, and we treat each image as Antisemitic/Islamophobic (depending on the phrase used for the comparison) if the cosine similarity between the phrase and the image is above the threshold. Then, we calculate the accuracy, precision, recall, and F1 score, for each of the ten phrases. We report the average performance across all phrases and the standard deviation (as the area) in Fig. 2. We observe that the model performs best with a cosine similarity threshold of 0.3 as we achieve a 0.84, 0.47, 0.80, 0.58 for accuracy, precision, recall, and f1 score, respectively. Indeed, the 0.3 threshold is also used by previous work by Schuhmann et al. [55] that inspected CLIP’s cosine similarities between text and images and determined that 0.3 is a suitable threshold.

Finally, to construct our Antisemitic/Islamophobic image dataset (*Visual* dataset), we extract all images that have a cosine similarity of 0.3 or higher with any of the Antisemitic/Islamophobic text phrases. We label each image as Antisemitic or Islamophobic depending on whether the textual phrase is Antisemitic or Islamophobic. To identify unique images, we use the Perceptual Hashing (pHash) algorithm [39] that calculates a fingerprint for each image in such a way that any two images that look similar to the human eye map have minor differences in their hashes. Similar to the Textual dataset, we remove all images labeled as both Antisemitic and Islamophobic (3,325 images), mainly because our manual inspections indicate that most of them are noise. Overall, we find 69,610 Antisemitic and 22,519 Islamophobic images that are shared in 472,048 and 101,465 posts, respectively (see Table 1).

Antisemitic phrases		Islamophobic phrases	
Phrase	# Posts	Phrase	# Posts
a kike	51,216	fuck muslim	7,993
fuck kike	23,604	kill muslim	4,639
fuck jew	20,241	fuck islam	3,464
gas the kike	13,353	kill all muslim	1,672
fuck off kike	11,108	muslim be terrorist	1,445
kike shill	10,134	i hate muslim	1,379
gas the kike race war now	6,105	muslim shithole	1,208
kill jew	5,815	muslim shit	1,085
you fuck kike	5,007	all muslim be terrorist	1,039
filthy kike	4,111	muslim be bad	1032
jew fuck	3,557	ban all muslim	958
kike faggot	3,540	fuck mudslimes	951
kike on a stick	3,537	muslim cunt	907
gas the jew	3,081	i hate islam	881
faggot kike	2,905	fuck sandniggers	876

**Table 2:** Top 15 phrases (lemmatized versions), in terms of the number of posts, in our Antisemitic and Islamophobic Textual dataset. For each phrase, we report the number of posts that contain it.



**Figure 4:** CDF of the number of Antisemitic/Islamophobic posts containing each image.

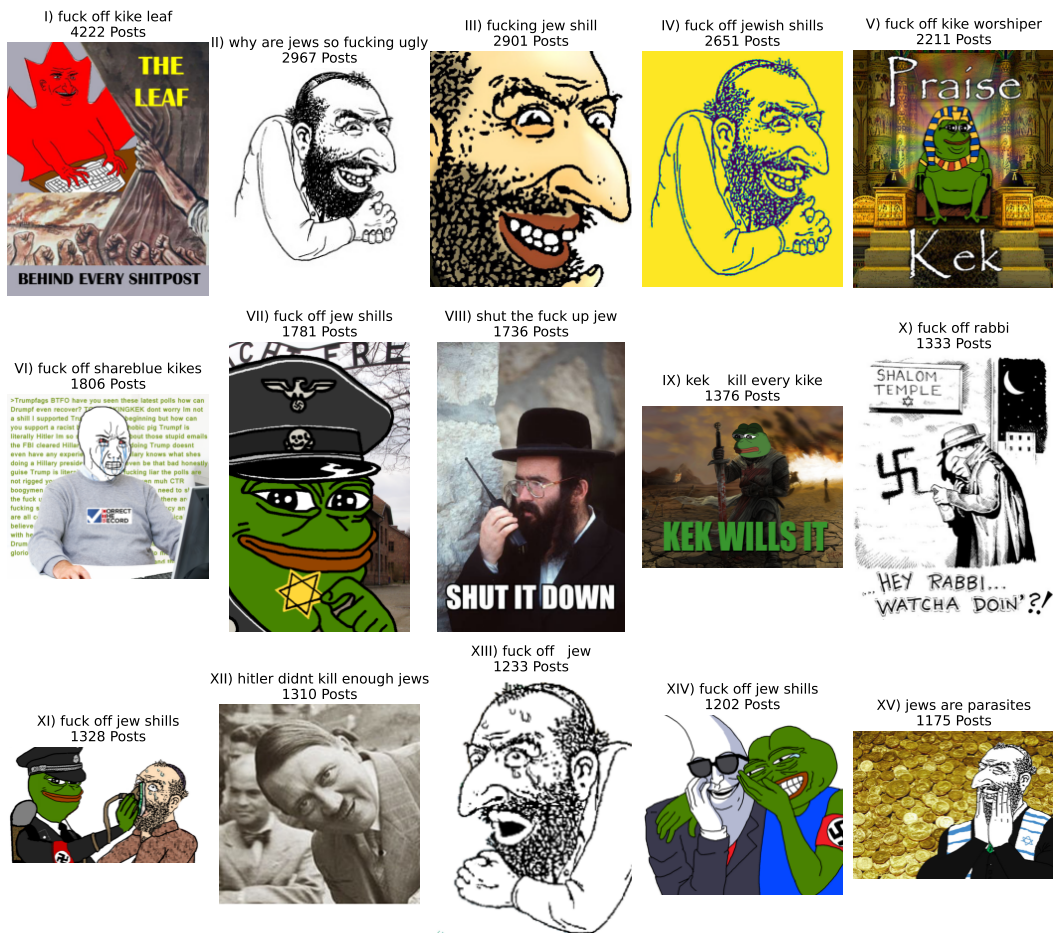
## 5 Results

This section presents our results from analyzing the Antisemitic/Islamophobic Textual and Visual datasets.

### 5.1 Popular phrases in Textual dataset

We start our analysis by looking into the most popular phrases in our Antisemitic/Islamophobic textual datasets. Fig. 3 shows the Cumulative Distribution Function (CDF) of the number of posts per each Antisemitic/Islamophobic phrase. We observe that these hateful phrases tend to appear in a considerable amount of posts. For instance, 90.4% and 92.47% of the Antisemitic and Islamophobic phrases appear in at least ten posts. Furthermore, we identify that the percentage of Antisemitic phrases (41.8%) that appear in at least 100 posts is slightly lower than the percentage of Islamophobic phrases (49.46%). At the same time, we observe that a small percentage of phrases (11.54%) is shared in more than 1000 posts on the Antisemitic/Islamophobic textual datasets combined.

We also report the top 15 phrases, in terms of the number of posts, in our Antisemitic and Islamophobic Textual dataset (see Table 2). In the first dataset, we observe that 12 out of the



**Figure 5:** Top 15 images, in terms of the number of posts, in our Antisemitic Visual datasets. We report the number of posts that the image appears and the most related Antisemitic phrase that matches each image.

15 most frequent phrases contain the term “kike”, a derogatory term to denote Jews. We also identify three phrases related to the extermination procedure in the gas chambers during the holocaust. Indeed, 16,433 (7.85%) of the Antisemitic posts contain at least one of these phrases: “gas the kike”, “gas the jew”, or “gas the kike race war now”. Phrases accusing jews of being accomplices (“kike shill”) or alluding to a supposed good social-economic status (“filthy kike”) are also trendy, appearing in 10,134 and 4,111 posts, respectively.

We also show the top 15 most popular Islamophobic phrases in Table 2. Here, we observe many posts with phrases calling terrorists to Muslims. For instance, “Muslims be terrorist” and “All Muslim be terrorist” appear in 1,445 and 1,039 posts, respectively. The second and fourth most popular phrases are calls for attacks targeting Muslims; “Kill Muslim” and “Kill all Muslim” appear in approximately 4.6K and 1.6K posts. We also find phrases against Islam; “Fuck Islam” appears in 3.4K posts and “I hate Islam” in 881 posts. Finally, we also identify phrases containing the terms “mudslimes” [59] and “sandniggers” [60], which are derogatory names to refer to Muslims and Arabs.

## 5.2 Popular Images in Visual Dataset

We also look into the popularity of images in our Antisemitic/Islamophobic datasets (in terms of the number of posts they shared). Fig. 4 shows the CDF of the number of posts for each Antisemitic/Islamophobic image. We observe that Antisemitism and Islamophobia imagery is a diverse problem, with 58.64% and 62.58% of the images appearing only in one post for Antisemitism and Islamophobia, respectively. At the same time, we have a small percentage of images that are shared many times on 4chan’s /pol/; 8.94% of all the Antisemitic/Islamophobic imagery are shared at least ten times. Overall, we observe a similar pattern in the distribution of the number of posts per image for Antisemitism and Islamophobia.

Next, we look into the most popular images in our Antisemitic and Islamophobic visual datasets. Fig. 5 shows the top 15 Antisemitic images; for each image, we report the number of posts it appears and the most relevant matching textual phrase. We observe that the Happy Merchant meme appears in a lot of the most popular Antisemitic memes on 4chan’s /pol/ (images I, II, III, IV, X, XI, XIII, and XV in Fig. 5). The most popular image is particularly interesting as it likely aims to disseminate the idea that Canadians are behind every shitpost on /pol/ (note the antisemitic connotation in the leaf as it has



**Figure 6:** Top 15 images, in terms of the number of posts, in our Islamophobic Visual dataset. We report the number of posts that the image appeared and the most related Islamophobic phrase for each image.

the Happy Merchant’s face). We also observe some false positives among the most popular images (image V, VI, and IX in Fig. 5). These images are not Antisemitic; however, by looking into the matching phrase, we can understand why the CLIP model treats it as similar. For instance, image V in Fig. 5 is likely returned because of the word “worshiper” in the query. This specific example highlights the sensitivity of the CLIP model to the input queries and how it can result in the generation of false positives when considering the task of hateful content detection. Other interesting examples of popular images include image VIII and image X in Fig. 5; both hinting that members of the Jewish are allegedly the masterminds of lousy stuff happening or conspiracy theories (i.e., shut down the Jewish plan or Rabbi painting Nazi symbols).

We also show the top 15 most popular Islamophobic images in Fig. 6. Here, we observe some images that are pretty graphic; for example, image V shows a pig having sex with a Muslim, or image VI shows a Muslim eating his excreta. Other popular images include ones that show Pepe the frog as a Muslim (these images are not directly Islamophobic, but CLIP considers them as related to the query), images that include sarcasm and link Muslims to terrorism (image VII and XIII in Fig. 6), as well as images linking Muslims to the Happy Merchant meme (e.g., image XIV in Fig. 6). Again, similarly to the most popular images in the Antisemitic dataset, we have

some false positives like image III (likely because the CLIP model does not know the slur “sandniggers”) and image XV in Fig. 6. Finally, the second most popular image in our Islamophobic dataset (Fig. 6) likely highlights that the CLIP model has some biases; e.g., it links the phrase “all muslims must hang” with an image showing Pepe the Frog dressed as a terrorist, likely indicating that the model thinks that Muslims are terrorists.

### 5.3 Antisemitic/Islamophobic content over time

This section presents our temporal analysis that shows the distribution of Antisemitic/Islamophobic content over time. Fig. 7 shows the number of hateful posts per day in the Antisemitic Textual/Visual datasets. We run Kendall’s tau-b correlation to determine the relationship between the number of posts in the Antisemitic Textual and Visual dataset. We find a strong, positive, and statistically significant correlation ( $\tau = .582, p < .001$ ), indicating that hateful content is spread both using text and images in a similar fashion. We also observe the highest volume of textual and image content between April 6, 2017, and April 9, 2017, with 4,132 (1.97% of the dataset) posts in the Textual dataset and 7,134 (1.51%) posts in the Visual dataset. This finding confirms previous findings from Zannettou et al. [71] that identified a spike in the spread

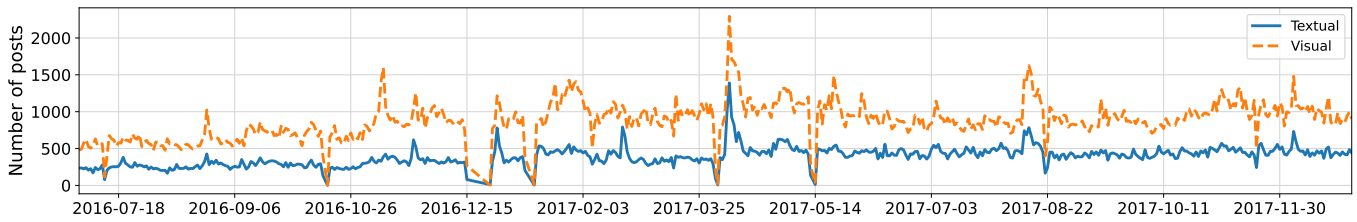


Figure 7: Number of Antisemitic posts per day in our Textual/Visual datasets.

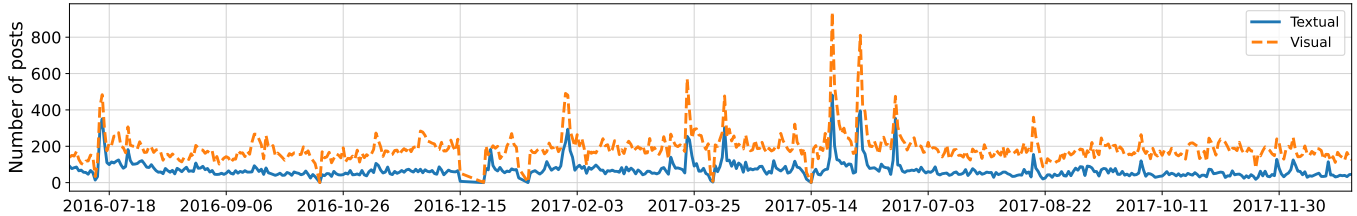


Figure 8: Number of Islamophobic posts per day in our Textual/Visual datasets.

of the Happy Merchant memes on April 7, 2017.

By inspecting the top 15 most frequent images during that period (we omit the figure due to space constraints), we identified that those images are related to the decision of Donald Trump to remove Steve Bannon from the National Security Council Post on April 5, 2017 [15] and a missile attack in Syria on April 7, 2017 [53]. According to newspapers [5, 27], Jared Kushner, the Jewish Trump’s son-in-law, seemed to be acting as a shadow secretary of state visiting and taking Middle East portfolios after that event. This political decision spread a volume of image content with the face of Jared Kushner. Also, there are some references to Donald Trump that indicate that he is controlled by Israel ( e.g., the most popular image is associated with the phrase “fuck Drumpf and fuck Jewish people”). The term “Drumpf” is a German surname most commonly known as the likely predecessor to the family name of Donald Trump.

We also evaluate the distribution of Islamophobic posts over time. Fig. 8 shows the number of Islamophobic posts per day in our Textual/Visual datasets. We also find a statistically significant, strong, and positive correlation ( $\tau = .470, p < .001$ ). In both datasets, we find a peak of activity on May 23, 2017, with 482 and 937 posts in the Textual and Visual datasets, respectively. By manually inspecting the top 15 images shared that day, we identify that the high volume of posts is related to the Manchester Bombing; on May 22, 2017, a British man detonated a suicide bomb in the foyer of the Manchester Arena as people were leaving a concert by pop singer Ariana Grande. On May 23, ISIS claimed responsibility for the attack. [14]. This event raised hateful online narratives defining Muslims as terrorists [20]. We find images that contain explicit references to this attack and images questioning whether Islam is a religion of peace. Overall, our findings highlight that both textual and visual hateful content is likely influenced by real-world events, with peaks of hateful activity observed during important real-world events that are related to the demographic groups we study.

## 6 Related Work

**Hate Speech detection.** Hate speech has recently received much research attention, with several works focusing on detecting hate speech in online social media. Initial research on hate speech analysis is typically oriented towards monolingual and single classification tasks due to the complexity of the task. They used simple methods such as dictionary lookup [26], bag of words [26], or SVM classifiers [37, 57]. Recent efforts are proposing multilingual and multitask learning by using deep learning models [24, 43, 62, 63]. While previous approaches to characterize and identify hate speech focus purely on the *content* posted in social media, some research efforts shift the focus towards detecting hateful users by exploiting other contextual data [1, 11, 49, 64]. Furthermore, other research efforts investigate to what extent the models trained to detect general abusive language generalize between different datasets labeled with different abusive language types [32, 38, 40, 52, 54]. While less explored, some work focus on multimodal settings, formed by text and images [16, 33]. Gomez et al. [25] build a large dataset for multimodal hate speech detection retrieved from Twitter using specific hateful seed keywords, finding that multimodal models do not outperform the unimodal text ones.

**Antisemitism.** Antisemitism have grown and proliferated rapidly online and have done so mostly unchecked; Zannettou et al. [71] call for new techniques to understand it better and combat it. Ozalp et al. [44] train a scalable supervised machine learning classifier to identify antisemitic content on Twitter. Chandra et al. [9] propose a multimodal system that uses text, images, and OCR to detect the presence of Antisemitic textual and visual content. They apply their model on Twitter and Gab, finding that multiple screenshots, multi-column text, and texts expressing irony, sarcasm posed problems for the classifiers. To characterize antisemitism, Enstad [21] propose an analytical framework composed of three indicators: antisemitic attitudes, incidents targeting Jews, and Jew’s exposure to antisemitism. Their results show that attitudes vary by geographic and cultural region and among population sub-groups.



**Islamophobia.** Surveys show that Islamophobia is rising on Web communities [29]. Vidgen and Yasseri [61] build an SVM classifier to distinguish between tweets non-Islamophobic, weak Islamophobic, and strong Islamophobic with a balanced accuracy of 83%. Cervi [8] use clause-based semantic text analysis to identify the presence of Islamophobia in electoral discourses of political parties from Spain and Italy. Chandra et al. [10] apply topic modeling and temporal analysis over tweets from the #coronajihad to identify the existence of Islamophobic rhetoric around COVID-19 in India. Civila et al. [13] apply content analysis over 474 images and texts in from Instagram posts under the hashtag #StopIslam. Alietti and Padovan [2] conduct telephone surveys on 1.5K Italians on Antisemitic and Islamophobic attitudes, finding an overlap of ideology for both types of hate speech.

## 7 Discussion & Conclusion

In this work, we explored the problem of Antisemitism/Islamophobia on 4chan’s /pol/ using OpenAI’s CLIP model. We devised a methodology to identify Antisemitic/Islamophobic textual phrases using Google’s Perspective API and manual annotations and then used the CLIP model to identify hateful imagery based on the phrases. We found that the CLIP can play a role in detecting hateful content; using our methods, the CLIP can detect hateful content with an accuracy of 84%. Also, we found that Antisemitic/Islamophobic imagery exists in 2x more posts when compared to Antisemitic/Islamophobic speech on 4chan’s /pol/. Additionally, our work contributes to research efforts focusing on detecting hateful content by making a dataset of 420 Antisemitic/Islamophobic phrases, 246K textual posts, and 92K images publicly available. Below, we discuss the implications of our findings for researchers focusing on detecting hate speech and for researchers working on large pre-trained models like OpenAI’s CLIP.

**Prevalence of Antisemitic/Islamophobic Imagery.** Our findings show that images play a significant role in the spread of hateful content, and on 4chan’s /pol/ they even overshadow hateful content spread via text. This is likely because 4chan is an imageboard and a fringe Web community; hence, a large volume of hateful content is disseminated via images. Nevertheless, the problem of hateful imagery exists on other mainstream platforms (e.g., Twitter), hence it is of paramount importance to develop better and more accurate systems for the detection of hateful content across multiple modalities. For instance, we argue that the spread of hateful content via videos is an unexplored problem, and there is a need to develop models across text, images, and videos.

**Performance and Sensitivity of CLIP model.** Our experiments indicate that large-pre-trained models like CLIP are pretty powerful and have general knowledge that can be used for various tasks. When considering the hateful content detection task, the CLIP model should be used with caution. This is because the CLIP model highly depends on how the input text query is written, influencing the number of false positives returned. When CLIP is used for moderation purposes, we ar-

gue that it is essential to have humans in the loop to ensure that the automated model works as expected. Additionally, we observed that the CLIP model performs worse when considering input text queries that comprise many words. This indicates that we need more powerful text encoders that can capture the primary meaning of textual phrases, irrespectively of how long they are.

**Biases on CLIP model.** Large pre-trained models like OpenAI’s CLIP are trained on large-scale datasets from the Web, and these datasets might include biases, hence some of the bias is transferred to the trained model. From our experiments and manual annotations, we observed some instances of such biases; e.g., the CLIP model identifying an image showing a terrorist as similar to a text phrase talking about Muslims (i.e., the model is biased towards Muslims, thinking they are terrorists). When considering that these models can be used for moderation purposes (e.g., detecting and removing hateful content), such biases can result in false positives biased towards specific demographics. This can cause users to lose trust in the platform and its moderation systems and may cause them to stop using the platform. Overall, given the increasing use of such models in real-world applications, there is a pressing need to develop techniques and tools to diminish such biases from large pre-trained models.

**Limitations.** Our work has several limitations. First, we rely on Google’s Perspective API to initially identify hateful text, which has its limitations (e.g., might not understand specific slurs posted on 4chan) and biases when detecting hateful text. Second, our analysis focuses only on short textual phrases (at most seven words), mainly because our preliminary results showed that CLIP does not perform well in detecting hateful imagery when considering long phrases. Therefore, it is likely to miss some Antisemitic/Islamophobic text and imagery. Third, we rely entirely on a pre-trained CLIP model; this is not ideal since the CLIP model is trained on a public dataset obtained from multiple Web resources and is not specific to our platform of interest (i.e., 4chan). This might result in the model not recognizing 4chan slurs or slang language. As part of our future work, we intend to explore the possibility of fine-tuning the CLIP model with datasets obtained from fringe Web communities like 4chan.

## References

- [1] Zo Ahmed, Bertie Vidgen, and Scott A. Hale. Tackling racial bias in automated online hate detection: Towards fair and accurate classification of hateful online users using geometric deep learning. *CoRR*, abs/2103.11806, 2021.
- [2] Alfredo Alietti and Dario Padovan. Religious racism, islamophobia and antisemitism in italian society. *Religions*, 4(4):584–602, 2013.
- [3] Anonymous. Common antisemitic and islamophobic phrases. <https://docs.google.com/spreadsheets/d/1zguFPsiU8zHJd9E9bka8cagsBBCVNCM4ofLuGvgXG5w/edit?usp=sharing>, 2022.
- [4] Aymé Arango, Jorge Pérez, and Barbara Poblete. Hate speech

- detection is not as easy as you may think: A closer look at model validation. In *SIGIR*, 2019.
- [5] Peter Baker, Maggie Haberman, and Glenn Thrush. Trump removes stephen bannon from national security council post. <https://www.nytimes.com/2017/04/05/us/politics/national-security-council-stephen-bannon.html>, 2017.
  - [6] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
  - [7] Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of supervised learning algorithms. In *ICML*, 2006.
  - [8] Laura Cervi. Exclusionary populism and islamophobia: A comparative analysis of italy and spain. *Religions*, 11(10):516, 2020.
  - [9] Mohit Chandra, Dheeraj Pailla, Himanshu Bhatia, Aadilmehdi Sanchawala, Manish Gupta, Manish Shrivastava, and Ponnurangam Kumaraguru. “subverting the jewtocracy”: Online anti-semitism detection using multimodal deep learning. In *WebSci*, 2021.
  - [10] Mohit Chandra, Manvith Reddy, Shradha Sehgal, Saurabh Gupta, Arun Balaji Buduru, and Ponnurangam Kumaraguru. “a virus has no religion”: Analyzing islamophobia on twitter during the covid-19 outbreak. In *HT*, 2021.
  - [11] Prateek Chaudhry and Matthew Lease. You are what you tweet: Profiling users by past tweets to improve hate speech detection. *arXiv preprint arXiv:2012.09090*, 2020.
  - [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.
  - [13] Sabina Civila, Luis M Romero-Rodríguez, and Amparo Civila. The demonization of islam through social media: A case study of #stopislam in instagram. *Publications*, 8(4):52, 2020.
  - [14] Ian Cobain, Frances Perraudin, Steven Morris, and Nazia Parveen. Salman ramadan abedi named by police as manchester arena attacker:”. *The Guardian*, 2017.
  - [15] Robert Costa and Abby Phillip. Stephen bannon removed from national security council. <https://www.washingtonpost.com/news/post-politics/wp/2017/04/05/stephen-bannon-no-longer-a-member-of-national-security-council/>, 2017.
  - [16] Abhishek Das, Japsimar Singh Wahi, and Siyao Li. Detecting hate speech in multi-modal memes. *arXiv preprint arXiv:2012.14891*, 2020.
  - [17] Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language. In *ICWSM*, 2017.
  - [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
  - [19] Ali Diba, Vivek Sharma, Reza Safdari, Dariush Lotfi, Saquib Sarfraz, Rainer Stiefelhofen, and Luc Van Gool. Vi2clr: Video and image for visual contrastive learning of representation. In *ICCV*, pages 1502–1512, 2021.
  - [20] Joseph Downing, Sarah Gerwens, and Richard Dron. Tweeting terrorism: Vernacular conceptions of muslims and terror in the wake of the manchester bombing on twitter. *CTS*, pages 1–28, 2022.
  - [21] Johannes D Enstad. Contemporary antisemitism in three dimensions: A new framework for analysis. *SocArXiv*, 28, 2021.
  - [22] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
  - [23] John M Giorgi, Osvald Nitski, Gary D Bader, and Bo Wang. Declutr: Deep contrastive learning for unsupervised textual representations. *arXiv preprint arXiv:2006.03659*, 2020.
  - [24] Goran Glavaš, Mladen Karan, and Ivan Vulić. XHate-999: Analyzing and detecting abusive language across domains and languages. In *COLING*, December 2020.
  - [25] Raul Gomez, Jaume Gibert, Lluís Gomez, and Dimosthenis Karatzas. Exploring hate speech detection in multimodal publications. In *WACV*, pages 1470–1478, 2020.
  - [26] Radhouane Guermazi, Mohamed Hammami, and Abdelmajid Ben Hamadou. Using a semi-automatic keyword dictionary for improving violent web site filtering. In *SITIS 2007*, pages 337–344, 2007.
  - [27] Maggie Haberman, Jeremy W Peters, and Peter Baker. In battle for trump’s heart and mind, it’s bannon vs. kushner. <https://www.nytimes.com/2017/04/06/us/politics/stephen-bannon-white-house.html>, 2017.
  - [28] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, 2006.
  - [29] Farid Hafez, Enes Bayrakli, Léonard Faytre, Amina Easat-Daas, Anna-Esther Younes, Natalia Kutuzova, Hikmet Karčić, Hayri A Emin, Nejra Kadić Meškić, Selma Muhic Dizdarevic, et al. European islamophobia report 2019. [https://setav.org/en/assets/uploads/2020/06/EIR\\_2019.pdf](https://setav.org/en/assets/uploads/2020/06/EIR_2019.pdf), 2019.
  - [30] Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive multi-view representation learning on graphs. In *ICML*, 2020.
  - [31] Gabriel Hine, Jeremiah Onalapo, Emiliano De Cristofaro, Nicolas Kourtellis, Ilias Leontiadis, Riginos Samaras, Gianluca Stringhini, and Jeremy Blackburn. Kek, cucks, and god emperor trump: A measurement study of 4chan’s politically incorrect forum and its effects on the web. In *ICWSM*, 2017.
  - [32] Mladen Karan and Jan Šnajder. Cross-domain detection of abusive language online. In *ALW2*, pages 132–137, Brussels, Belgium, October 2018. Association for Computational Linguistics.
  - [33] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. *arXiv preprint arXiv:2005.04790*, 2020.
  - [34] Sungnyun Kim, Gihun Lee, Sangmin Bae, and Se-Young Yun. Mixco: Mix-up contrastive learning for visual representation. *arXiv preprint arXiv:2010.06300*, 2020.
  - [35] Daniel Konikoff. Gatekeepers of toxicity: Reconceptualizing twitter’s abuse and hate speech policies. *P&I*, 2021.
  - [36] Tarald O. Kvalseth. Note on cohen’s kappa. *Psychol. Rep.*, 1989.
  - [37] Shervin Malmasi and Marcos Zampieri. Detecting hate speech in social media. In *RANLP*, 2017.
  - [38] Johannes Skjeggstad Meyer and Björn Gambäck. A platform agnostic dual-strand hate speech detector. In *ALW*, pages 146–156, Florence, Italy, August 2019. Association for Computational Linguistics.
  - [39] Vishal Monga and Brian L. Evans. Perceptual image hashing via feature points: Performance evaluation and tradeoffs. *TIP*, 15:3452–3465, 2006.

- [40] Isar Nejadgholi and Svetlana Kiritchenko. On cross-dataset generalization in automatic detection of online abuse. *arXiv preprint arXiv:2010.07414*, 2020.
- [41] NLTK. Nltk tokenization. <https://www.nltk.org/api/nltk.tokenize.html>, 2021.
- [42] NLTK. Nltk lemmatization. [https://www.nltk.org/\\_modules/nltk/stem/wordnet.html](https://www.nltk.org/_modules/nltk/stem/wordnet.html), 2021.
- [43] Nedjma Ousidhoum, Zizheng Lin, Hongming Zhang, Yangqiu Song, and Dit-Yan Yeung. Multilingual and multi-aspect hate speech analysis. In *EMNLP-IJCNLP*, pages 4675–4684, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [44] Sefa Ozalp, Matthew L. Williams, Pete Burnap, Han Liu, and Mohamed Mostafa. Antisemitism on twitter: Collective efficacy and the role of community organisations in challenging online hate speech. *Soc. Media Soc.*, 2020.
- [45] Antonis Pappasavva, Savvas Zannettou, Emiliano De Cristofaro, Gianluca Stringhini, and Jeremy Blackburn. Raiders of the lost kek: 3.5 years of augmented 4chan posts from the politically incorrect board. In *ICWSM*, 2020.
- [46] Perspective API. <https://www.perspectiveapi.com/>, 2018.
- [47] Dan Prisk. The hyperreality of the alt right: how meme magic works to create a space for far right politics. 2017.
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [49] Manoel Horta Ribeiro, Pedro H Calais, Yuri A Santos, Virgílio AF Almeida, and Wagner Meira Jr. Characterizing and detecting hateful users on twitter. In *ICWSM*, 2018.
- [50] Manoel Horta Ribeiro, Shagun Jhaver, Savvas Zannettou, Jeremy Blackburn, Emiliano De Cristofaro, Gianluca Stringhini, and Robert West. Does platform migration compromise content moderation? evidence from r/the\_donald and r/incels. In *CSCW*, 2021.
- [51] Caitlin M Rivers and Bryan L Lewis. Ethical research standards in a world of big data. *F1000Research*, 3, 2014.
- [52] Marian-Andrei Rizoiu, Tianyu Wang, Gabriela Ferraro, and Hanna Suominen. Transfer learning for hate speech detection in social media. *arXiv preprint arXiv:1906.03829*, 2019.
- [53] Everett Rosenfeld. Trump launches attack on syria with 59 tomahawk missiles. <https://www.cnn.com/2017/04/06/us-military-has-launched-more-50-than-missiles-aimed-at-syria-nbc-news.html>, 2017.
- [54] Joni Salminen, Maximilian Hopf, Shammur A Chowdhury, Soon-gyo Jung, Hind Almerkhi, and Bernard J Jansen. Developing an online hate classifier for multiple social media platforms. *HCIS*, 10(1):1–34, 2020.
- [55] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- [56] Andrew Sellars. Defining hate speech. *Berkman Klein Center Research Publication*, (2016-20):16–48, 2016.
- [57] Yonas Senarath and Hemant Purohit. Evaluating semantic feature representations to efficiently detect hate intent on social media. In *2020 IEEE 14th International Conference on Semantic Computing (ICSC)*, pages 199–202, 2020.
- [58] Marc Tuters, Emilija Jokubauskaitė, and Daniel Bach. Post-truth protest: how 4chan cooked up the pizzagate bullshit. *M/c Journal*, 21(3), 2018.
- [59] Urban Dictionary. Mudslime definition. <https://www.urbandictionary.com/define.php?term=mudslime>, 2006.
- [60] Urban Dictionary. Sandniggers definition. <https://www.urbandictionary.com/define.php?term=sand%20niggers>, 2006.
- [61] Bertie Vidgen and Taha Yasseri. Detecting weak and strong islamophobic hate speech on social media. *J. Inf. Technol. Politics*, 17(1):66–78, 2020.
- [62] Fedor Vitiugin, Yonas Senarath, and Hemant Purohit. Efficient detection of multilingual hate speech by using interactive attention network with minimal human feedback. In *WebSci*, 2021.
- [63] Kunze Wang, Dong Lu, Caren Han, Siyu Long, and Josiah Poon. Detect all abuse! toward universal abusive language detection models. In *COLING*, pages 6366–6376, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics.
- [64] Zeerak Waseem and Dirk Hovy. Hateful symbols or hateful people? predictive features for hate speech detection on Twitter. In *NAACL*, pages 88–93, San Diego, California, June 2016. Association for Computational Linguistics.
- [65] Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madihan Khabsa, Fei Sun, and Hao Ma. Clear: Contrastive learning for sentence representation. *arXiv preprint arXiv:2012.15466*, 2020.
- [66] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. *Adv. Neural Inf. Process. Syst.*, 33:5812–5823, 2020.
- [67] Xin Yuan, Zhe Lin, Jason Kuen, Jianming Zhang, Yilin Wang, Michael Maire, Ajinkya Kale, and Baldo Faieta. Multimodal contrastive training for visual representation learning. In *CVPR*, 2021.
- [68] Savvas Zannettou, Tristan Caulfield, Emiliano De Cristofaro, Nicolas Kourtellis, Ilias Leontiadis, Michael Sirivianos, Gianluca Stringhini, and Jeremy Blackburn. The web centipede: understanding how web communities influence each other through the lens of mainstream and alternative news sources. In *IMC*, 2017.
- [69] Savvas Zannettou, Tristan Caulfield, Jeremy Blackburn, Emiliano De Cristofaro, Michael Sirivianos, Gianluca Stringhini, and Guillermo Suarez-Tangil. On the origins of memes by means of fringe web communities. In *IMC*, 2018.
- [70] Savvas Zannettou, Mai ElSherief, Elizabeth Belding, Shirin Nilizadeh, and Gianluca Stringhini. Measuring and characterizing hate speech on news websites. In *WebSci*, 2020.
- [71] Savvas Zannettou, Joel Finkelstein, Barry Bradlyn, and Jeremy Blackburn. A quantitative approach to understanding online antisemitism. In *ICWSM*, 2020.
- [72] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020.