# **Detecting Abusive Language on Online Platforms: A Critical Analysis**

 $\begin{array}{c} \textbf{Preslav Nakov}^{1,2*}\,,\,\, \textbf{Vibha Nayak}^1\,,\,\, \textbf{Kyle Dent}^1\,,\,\, \textbf{Ameya Bhatawdekar}^3\\ \textbf{Sheikh Muhammad Sarwar}^{1,4}\,,\,\, \textbf{Momchil Hardalov}^{1,5},\, \textbf{Yoan Dinkov}^1\\ \textbf{Dimitrina Zlatkova}^1\,,\,\, \textbf{Guillaume Bouchard}^1\,,\,\, \textbf{Isabelle Augenstein}^{1,6} \end{array}$ 

<sup>1</sup>CheckStep Ltd., <sup>2</sup>Qatar Computing Research Institute, HBKU, <sup>3</sup>Microsoft, <sup>4</sup>University of Massachusetts, Amherst, <sup>5</sup>Sofia University, <sup>6</sup>University of Copenhagen {preslav.nakov, vibha, kyle.dent, momchil, yoan.dinkov, didi, guillaume, isabelle} @checkstep.com, ambhataw@microsoft.com, smsarwar@cs.umass.edu,

#### **Abstract**

Abusive language on online platforms is a major societal problem, often leading to important societal problems such as the marginalisation of underrepresented minorities. There are many different forms of abusive language such as hate speech, profanity, and cyber-bullying, and online platforms seek to moderate it in order to limit societal harm, to comply with legislation, and to create a more inclusive environment for their users. Within the field of Natural Language Processing, researchers have developed different methods for automatically detecting abusive language, often focusing on specific subproblems or on narrow communities, as what is considered abusive language very much differs by context. We argue that there is currently a dichotomy between what types of abusive language online platforms seek to curb, and what research efforts there are to automatically detect abusive language. We thus survey existing methods as well as content moderation policies by online platforms in this light, and we suggest directions for future

# 1 Introduction

Online harm is not new. Groups and individuals who have been targets of abusive language have suffered real harms for many years. The problem has persisted over time and may well be growing. Thus, abusive language in social media is of particular concern to online communities, governments, and especially social media platforms.

While combating abuse is a high priority, preserving individuals' rights to free expression is also vital, making the task of moderating conversations particularly difficult. Some form of moderation is clearly required. Platform providers face very challenging technical and logistical problems in limiting abusive language, while at the same time allowing a level of free speech which leads to rich and productive online conversations. Negative interactions experienced by users pose a significant risk to these social platforms and can adversely

affect not only user engagement, but can also erode trust in the platform and hurt a company's brand.

Social platforms have to strike the right balance in terms of managing a community where users feel empowered to engage while taking steps to effectively mitigate negative experiences. They need to ensure that their users feel safe, their personal privacy and information is protected, and that they do not experience harassment or annoyances, while at the same time feeling empowered to share information, experiences, and views. Many social platforms institute guidelines and policies to specify what content is considered inappropriate. As manual filtering is hard to scale, and can even cause post-traumatic stress disorder-like symptoms to human annotators, there have been many research efforts to develop tools and technology to automate some of this effort.

The key feature of offensive language in online dialog is that it is harmful either to its target, the online community where it occurs, or the platform hosting the conversation. The degree of harm is a factor and might range from hate speech and cyber-bullying with extreme deleterious effects, over slightly less damaging derogatory language and personal insults, to profanity and teasing, which might even be considered acceptable in some communities. This spectrum poses challenges for clear labeling of training data, as well as for computational modeling of the problem.

Several studies have considered the application of computational methods to deal with offensive language, particularly for English [Davidson et al., 2017; Basile et al., 2019; Fortuna and Nunes, 2018]. ever, these efforts tend to constrain their scope to one medium, to a single or to just a few subtasks, or to only limited aspects of the problem. Prior work has studied offensive language on Twitter [Xu et al., 2012; Burnap and Williams, 2015; Davidson et al., 2017; Wiegand et al., 2018], in Wikipedia comments, and in Facebook posts [Kumar et al., 2018]. The task is usually modeled as a supervised classification problem. The models are trained on posts annotated according to the presence of some type of abusive or offensive content. Examples of such types content include hate speech [Davidson et al., 2017; Malmasi and Zampieri, 2017;

<sup>\*</sup>Contact Author

<sup>&</sup>lt;sup>1</sup>challengehttps://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge

Platform Type	Example
Social Media	Twitter, Facebook, Instagram
Online Marketplace	Amazon, PinDuoDuo, Depop
Dating	Bumble, Tinder, Hinge
Video Community	TikTok, Triller, Clash
Forum	Patient, BG Mamma, Reddit
Gaming Community	Twitch, DLive, Omlet Arcade

Table 1: Classification of online platforms.

Malmasi and Zampieri, 2018], abusive language [Founta et al., 2018], toxicity [Georgakopoulos et al., 2018], cyber-bullying [Dinakar et al., 2011], and aggression [Kumar et al., 2018].

While there have been several surveys on offensive language, hate speech, and cyber-bullying, none of them have focused on what platforms need vs. what technology has to offer. For example, [Schmidt and Wiegand, 2017] and [Fortuna and Nunes, 2018] both surveyed automated hate speech detection, but focused primarily on the features which have been shown to be most effective in classification systems. [Salawu et al., 2020] provided an extensive survey of NLP for detecting cyber-bullying, and [Vidgen and Derczynski, 2021] did important work to catalog abusive language training data, which serves as a pre-step for system solutions. Here, we aim to bridge the gap between this work and platform solution requirements.

# 2 Requirements by Online Platforms

Below, we explore how abusive language differs from one online platform to another, and what content moderation policies have been put in place accordingly.<sup>2</sup>

Online platforms is a broad term which represents various categories of providers such as social media, online market-places, online video communities, dating websites and apps, support communities or forums, and online gaming communities, among others. Each of these platforms is governed by their own content moderation policies and has its own definitions of what constitutes abusive language.

Table 1 shows a classification of online platforms by type. Facebook, Instagram and Twitter represent big social media platforms, while Amazon, TikTok, Patient and Twitch are examples of platforms focusing on specific products.

### 2.1 The Big Tech

The big social media platforms, also commonly known as 'Big Tech', have stringent content moderation policies and the most advanced technology to help detect abusive language. We summarize these policies in Table 2.

We can see a lot of overlap, but also some differences in these policies.

### For **Facebook**<sup>3</sup>:

- Misinformation is always referenced as false news in Facebook's Community Guidelines.
- Even though Facebook covers most areas which could be subjected to abusive language, the coverage of *medical advice* is very broadly mentioned in its policies. It is unclear whether individuals are free to share medical advice to others in their posts.

### For **Twitter**<sup>4</sup>:

- The term *animal abuse* is not explicitly mentioned in the policy. It is implicitly mentioned under Twitter's *Sensitive media policy*, where graphic and sexual content featuring animals is not prohibited.
- Hate speech is referred to as hateful conduct in Twitter's content policy.
- Only COVID-19 related medical advice is prohibited on Twitter.

### For **Google**<sup>5</sup>:

- Google's terms of service covers only basic guidelines on acceptable conduct on the platform. The more specific clauses against hate speech, bullying and harassment are covered in service-specific policies as indicated in Table 2.
- There are no specific clauses addressing incidents of revenge porn and sexual abuse (adults). It is also unclear whether they fall under the category of illegal activities.

Amazon and Apple offer very different services in comparison to Facebook, Twitter and Google, which is reflected in the clauses covered in their terms of service.

## For **Apple** $^6$ :

- Apple's policies very broadly mention clauses in response to Dangerous people, Glorifying crime, Illegal goods, Child sexual abuse, Sexual abuse (Adults), Animal abuse, Human trafficking under illegal acts. Violence falls under threatening and intimidating posts.
- There are no clauses in the policy which address Medical advice, Spam, Sexual solicitation, Revenge porn, Graphic content, and Self-Harm.

### For **Amazon**<sup>7</sup>:

- Dangerous organization and people, Glorifying crime, Illegal goods, Child sexual abuse, Sexual abuse (Adults), Animal abuse, Human trafficking are broadly mentioned in clauses pertaining to Illegal.
- Revenge porn, Graphic content, Nudity and pornography are broadly mentioned in clauses pertaining to obscenity.
- There are no policies in place which directly address Medical advice, Misinformation, Sexual solicitation, and Self-Harm.

<sup>&</sup>lt;sup>2</sup>We link to each of these content policies, and additionally provide a static version of these policies at the time of writing here: https://bit.ly/37QEFmh

<sup>&</sup>lt;sup>3</sup>https://www.facebook.com/communitystandards/

<sup>&</sup>lt;sup>4</sup>https://help.twitter.com/en/rules-and-policies/twitter-rules

<sup>&</sup>lt;sup>5</sup>https://policies.google.com/terms?hl=en

<sup>&</sup>lt;sup>6</sup>https://www.apple.com/legal/internet-services/terms/site.html; https://support.apple.com/en-us/HT208863

<sup>&</sup>lt;sup>7</sup>https://www.amazon.com/gp/help/customer/display.html/?nodeId=508088

Policy Clauses	Facebook <sup>‡</sup> Twitter		Google	Apple	Amazon	
Violence	~	<b>✓</b>	<b>✓</b>	Intimidating, Threatening	◆ Threatening	
Dangerous orgs/people	✓	✓	Maps, Gmail, Meet*	? Illegal act	? under illegal	
Glorifying crime	✓	✓	Maps, Gmail, Meet*	? Illegal act	? under illegal	
Illegal goods	✓	✓	Maps, Google Chat and Hangout, Drive, Meet*	? Illegal act	? under illegal	
Self-harm	✓	✓	<b>✓</b>	×	×	
Child sexual abuse	✓	✓	<b>✓</b>	?	×	
Sexual Abuse (Adults)	✓	✓	×	?	?	
Animal abuse	✓	? Sensitive media policy	Earth, Drive, Meet*	? Illegal act	? under illegal	
Human trafficking	✓	<b>✓</b>	<b>✓</b>	? Illegal act	? under illegal	
Bullying and harassment	✓	✓	<b>✓</b>	<b>~</b>	◆ Threatening	
Revenge porn	✓	✓	×	×	? obscene	
Hate Speech	✓	✓ Hateful Conduct	<b>✓</b>	✓	Threatening	
Graphic content	✓	✓	Maps*	×	?	
Nudity and pornography	✓	✓	Earth, Meet, Drive, Chat and Hangout*	✓	? under obscene	
Sexual Solicitation	✓	×	Maps*	×	×	
Spam	✓	✓	<b>⋰</b>	×	✓	
Impersonation	✓	✓	Maps, Earth, Chat and Hangout, Gmail, Meet*	✓	✓	
Misinformation	✓ False news	✓	Maps, Drive*	✓	×	
Medical Advice	?	COVID-19 specific	Drive*	×	×	

Table 2: Policy Clauses of Big Tech. ✓ – Mentioned in the policy; 🛪 – Not mentioned in the policy; � – Implicitly mentioned in the policy; ¬ – Broadly mentioned under a clause in the policy; ¬ additional service specific policies; ¬ – the same policies are applied in Instagram.

### 2.2 Product-Specific Platforms

Content moderation policies also differ for single product platforms such as dating websites and apps, forums, online gaming and video communities, as can be seen in Table 3.

Since **Bumble** and **TikTok** are two of the most popular platforms, their policies cover most of the Big Tech clauses; however, with gaming communities and video communities, some of the clauses are either implicitly stated or are broadly described under a main clause.

One criticism for Big Tech's terms of service is that they are not explicit enough in the main terms. If one were to build machine learning models to flag content based on these, they would likely fail to capture their implicit meanings and thereby allow such content to persist. Inconsistencies across platforms also make modeling difficult. For instance, it is unclear whether *political propaganda* is allowed on Bumble, but it is explicitly disallowed by Facebook's policies.

### 3 Automatic Abusive Language Detection

Given the variety of needs across the different types of platforms, NLP models must be comprehensive and need to provide flexibility to capture differing priorities across disparate guidelines. Prior work, described below, has addressed several aspects of the problem, but we are unaware of any work which covers it comprehensively.

[Vidgen and Derczynski, 2021] attempted to categorize the available datasets for abusive language detection. In their survey, they analyzed 64 datasets, creating an overview of the landscape for datasets constructed in the past five years. Around half the datasets cover the English language text only.

Where multilinguality exists, it is presented mostly by European languages. As for Asian languages, datasets exist for Hindi and Indonesian, and six datasets contain Arabic text. The primary source of the data is Twitter. While other social networks are presented as well, they have fewer datasets. The size of most of the datasets is under 50,000 instances, and under 5,000 instances for half the datasets.

### 3.1 Target-Agnostic Schemata and Tasks

There have been several efforts to detect specific types of offensive content, e.g., hate speech, offensive language, cyberbullying, and cyber-aggression. Below, we briefly describe some of these datasets and the corresponding approaches.

**Hate speech identification** This is by far the most studied abusive language detection task [Kwok and Wang, 2013; Djuric et al., 2015; Burnap and Williams, 2015; Ousidhoum et al., 2019; Chung et al., 2019]. One of the most widely used datasets is the one by [Davidson et al., 2017], which contains over 24,000 English tweets labeled as nonoffensive, hate speech, and profanity. A recent shared task on the topic is HateEval [Basile et al., 2019] for English and Spanish. The problem was also studied from a multimodal perspective, e.g., [Sabat et al., 2019] developed a collection of 5,020 memes for hate speech detection. More recently, the Hateful Memes Challenge by Facebook introduced a dataset consisting of more than 10K memes, annotated as hateful or non-hateful [Kiela et al., 2020]: the memes were generated artificially, such that they resemble real memes shared on social media, along with 'benign confounders.'

Offensive language identification There have been several shared tasks with associated datasets, which focused specifically on offensive language identification, often featuring multiple languages: OffensEval 2019-2020 [Zampieri et al., 2019b; Zampieri et al., 2020] for English, Arabic, Danish, Greek, and Turkish, GermEval 2018 [Wiegand et al., 2018] HASOC 2019 [Mandl et al., 2019] for for German, German, and Hindi, TRAC 2018-2020 for English, Hindi [Fortuna et al., 2018; English, and Bengali,

https://bumble.com/en/guidelines;https://bumble.com/en/terms

<sup>8</sup>https://www.tiktok.com/community-guidelines?lang=en

<sup>&</sup>lt;sup>9</sup>https://patient.info/terms-and-conditions

<sup>&</sup>lt;sup>10</sup>https://www.redditinc.com/policies/content-policy; https://www.reddithelp.com/hc/en-us/articles/205926439

<sup>&</sup>lt;sup>11</sup>https://www.twitch.tv/p/en/legal/terms-of-service/; https://www.twitch.tv/p/en/legal/community-guidelines/

Policy Clauses	Bumble <sup>7</sup>	TikTok <sup>8</sup>	Patient* 9	$\mathbf{Reddit}^{10}$	Twitch <sup>11</sup>
Violence	<b>✓</b>	~	<b>✓</b>	<b>✓</b>	<b>✓</b>
Dangerous orgs/people	✓	~	unlawful act	under illegal	✓
Glorifying crime	✓	~	unlawful act	violent content	✓
Illegal goods	✓	~	unlawful act	✓	illegal activity
Self-harm	×	~	×	suicide response	<b>~</b> .
Child sexual abuse	✓	~	unlawful act	<b>✓</b> `	✓
Sexual Abuse (Adults)	illegal activity	~	unlawful act	illegal activity	✓
Animal abuse	×	~	unlawful act	violent content	illegal activity
Human trafficking	illegal activity	~	unlawful activity	illegal activity	
Bullying and harassment	<b>~</b>	~	<b>✓</b>	<b>~</b>	<b>~</b>
Revenge porn	×	×	×	? fake nudity	✓
Hate Speech	✓	~	✓	<b>✓</b> *	✓
Graphic content	✓	~	×	✓	✓
Nudity and pornography	✓	~	✓	✓	✓
Sexual Solicitation	✓	~	✓	✓	✓
Spam	✓	~	×	✓	✓
Impersonation	✓	~	<b>✓</b>	✓	✓
Misinformation	×	~	×	✓	✓
Medical Advice	×	~	<b>✓</b>	×	×

Table 3: Policy of Product-specific platforms; \* Patient is a medical advice forum.

Kumar et al., 2020]. Offensive language was also studied from a multimodal perspective, e.g., [Mittos et al., 2020] developed a dataset of memes shared on 4chan, and manually labeled as *offensive* vs. *non-offensive*.

Aggression identification The TRAC shared task on Aggression Identification [Kumar et al., 2018] provided participants with a dataset containing 15,000 annotated Facebook posts and comments in English and Hindi for training and validation. For testing, two different sets, one from Facebook and one from Twitter, were used. The goal was to discriminate between three classes: non-aggressive, covertly aggressive, and overtly aggressive. The best-performing systems in this competition used deep learning approaches based on convolutional neural networks (CNN), recurrent neural networks, and LSTMs [Aroyehun and Gelbukh, 2018; Majumder et al., 2018].

Toxic comment detection The Toxic Comment Classification Challenge<sup>1</sup> was an open competition at Kaggle, which provided participants with almost 160K comments from Wikipedia organized in six classes: toxic, severe toxic, obscene, threat, insult, and identity hate. The dataset was also used outside of the competition [Georgakopoulos et al., 2018], including as additional training material for the aforementioned TRAC shared task [Fortuna et al., 2018]. The task was later extended to multiple languages<sup>12</sup>, offering 8,000 Italian, Spanish, and Turkish comments. Recently, [Juuti et al, 2020] presented a systematic study of data augmentation techniques in combination with state-of-the-art pre-trained Transformer models for toxic language. A related Kaggle challenge features Detecting Insults in Social Commentary. Other datasets include Wikipedia Detox [Wulczyn et al., 2017], and those from [Davidson et al., 2017] and [Wulczyn et al., 2017].

**Cyberbullying detection** There have been several studies on cyberbullying detection. For example, [Xu et al., 2012] used sentiment analysis and topic models to identify relevant topics, and [Dadvar et al., 2013] employed user-related features such as the frequency of profanity in previous messages. [Rosa et al, 2019] presented a systematic review of automatic cyberbullying detection.

Abusive language detection There have also been datasets that cover various types of abusive language. [Founta et al., 2018] tackled hate and abusive speech on Twitter, introducing a dataset of 100K tweets. [Glavaš et. al, 2020] targeted hate speech, aggression, and attacks in three different domains: Fox News (from GAO), Twitter/Facebook (from TRAC), and Wikipedia (from WUL). In addition to English, it further offered parallel examples in Albanian, Croatian, German, Russian, and Turkish. However, the dataset is small, containing only 999 examples. Among the popular approaches for abusive language detection are cross-lingual embeddings [Ranasinghe and Zampieri, 2020] and deep learning [Founta et al., 2019].

### 3.2 Target-Aware Schemata and Tasks

As abusive language has many aspects, there have been several proposals for multi-level taxonomies covering different types of abuse. This permits understanding how different types and targets of abusive language relate to each other, and informs efforts to detect them. For example, offensive messages directed at a group are likely hate speech, whereas similar language targeting an individual is likely cyberbullying.

[Waseem et al., 2017] presented a taxonomy that differentiates between (abusive) language directed towards a specific individual or entity, or towards a generalized group, and whether the abusive content is explicit or implicit. [Wiegand et al., 2018] extended this idea to German tweets. They developed a model to detect offensive vs. non-offensive tweets, and further sub-classified the offensive tweets as profanity, insult, or abuse.

<sup>&</sup>lt;sup>12</sup>classificationhttps://www.kaggle.com/c/jigsaw-multilingual-toxic-comment-classification

<sup>&</sup>lt;sup>13</sup>http://www.kaggle.com/c/detecting-insults-in-social-commentary

Tweet	A	В	С
@USER Anyone care what that dirtbag says?	OFF	TIN	IND
Poor sad liberals. No hope for them.	OFF	TIN	GRP
LMAOYOU SUCK NFL	OFF	TIN	OTH
@USER What insanely ridiculous bullshit.	OFF	UNT	-
@USER you are also the king of taste	NOT	-	-

Table 4: Examples from the *OLID* dataset.

[Zampieri et al., 2019a] developed a very popular threelevel taxonomy which considers both the type and the target of offense:

Level A *Offensive Language Detection* Is it offensive?

OFF Inappropriate language, insults, or threats.

NOT Neither offensive, nor profane.

**Level B** Categorization of Offensive Language Is the offensive text targeted?

**TIN** Targeted insult or threat towards a group or individual.

**UNT** Untargeted profanity or swearing.

**Level C** *Offensive Language Target Identification* What is the target of the offense?

**IND** The target is an individual explicitly or implicitly mentioned in the conversation;

**GRP** Hate speech, targeting a group of people based on ethnicity, gender, sexual orientation, religion, or other common characteristic.

**OTH** Targets that do not fall into the previous categories, e.g., organizations, events, and issues.

Examples are shown in Table 4. The taxonomy served as the basis of the *OLID* dataset of 14,000 English tweets, which was used in two shared tasks (OffensEval) at SemEval in 2019 [Zampieri et al., 2019b] and in 2020 [Zampieri et al., 2020]. For the latter, an additional large-scale dataset was developed, consisting of nine million English tweets labeled in a semi-supervised fashion [Rosenthal et al., 2020]. This new dataset enabled sizable performance gains, especially at the lower levels of the taxonomy. The taxonomy was also successfully adopted for languages such as Arabic [Mubarak et al., 2020], Danish [Sigurbergsson and Derczynski, 2020], Greek [Pitenis et al., 2020], and Turkish [Çöltekin, 2020].

### 3.3 User-Aware Perspective

**Troll detection** has been addressed using semantic analysis [Cambria et al., 2010], domain-adapted sentiment analysis [Seah et al., 2015], various lexico-syntactic features about user writing style and structure [Chen et al., 2012], as well as graph-based approaches [Kumar et al., 2014]. There have also been studies on general troll behavior [Herring et al., 2002; Buckels *et al.*, 2014], cyberbullying [Galán-García et al., 2014; Sarna and Bhatia, 2017; Wong et al., 2018; Sezer et al., 2015], as well as on linking fake troll profiles to real users [Galán-García et al., 2014].

Some studies related to cyberbullying have already been applied in real settings in order to detect and to stop cyberbullying in elementary schools using a supervised machine learning algorithm that links fake profiles to real ones on the same social media [Galán-García et al., 2014].

Identification of malicious accounts in social networks is another important research direction. This includes detecting spam accounts [Almaatouq et al., 2016; Mccord and Chuah, 2011], fake accounts [Fire et al., 2014; Cresci et al., 2015], compromised accounts and phishing accounts [Adewole et al., 2017]. Fake profile detection has also been studied in the context of cyberbullying [Galán-García et al., 2014]. A related problem is that of *Web spam detection*, which has been addressed as a text classification problem [Sebastiani, 2002], e.g., using spam keyword spotting [Dave et al., 2003], lexical affinity of arbitrary words to spam content [Hu and Liu, 2004], frequency of punctuation and word co-occurrence [Li et al., 2006].

### 3.4 Learning Approaches

Nowadays, the most common way to address the above tasks is to use pre-trained transformers: typically BERT, but also RoBERTa, ALBERT [Lan et al., 2019], and GPT-2 [Radford et al., 2019]. In a multi-lingual setup, also mBERT [Devlin et al., 2019] and XLM-RoBERTa [Conneau et al, 2020] have proved useful. Other popular models include CNNs [Fukushima, 1980], RNNs [Rumelhart et al., 1986], and GRUs [Cho et al., 2014], including ELMo [Peters et al., 2018]. Older models such as SVMs [Cortes and Vapnik, 1995] are sometimes also used, typically as part of ensembles. Moreover, lexica such as HurtLex [Bassignana et al., 2018] and Hatebase are sometimes used as well.

For hateful meme detection, popular approaches include Visual BERT, ViLBERT, VLP, UNITER, LXMERT, VILLA, ERNIE-Vil, Oscar and various Transformers [Li et al., 2019; Su et al., 2020; Zhou et al., 2019; Tan and Bansal, 2019; Gan et al., 2020; Yu et al., 2020; Li et al., 2020; Vaswani et al., 2017; Lippe et al., 2020; Zhu, 2020; Muennighoff, 2020; Zhang et al., 2020].

## 4 Lessons Learned and Major Challenges

Despite the stringent policies in place, there are increasing cases of hate speech on social media platforms. Different platforms have different definitions for hate speech. Each policy tries to imagine every possible outcome of particular statements and tries to be transparent about the areas they fall short in. The reality is that some clauses in the terms of service require human attention and understanding. The following describes these challenges in more detail.

**Mismatch in goals.** As is readily apparent when comparing content moderation policies on online platforms (Section 2) with abusive language tasks tackled in NLP (Section 3) as well as methods researched (Section 3.4), there is a dichotomy between what types of content platforms seek to moderate, and what current abusive language detection

<sup>14</sup>http://hatebase.org/

method research offers. Most crucially, platforms need to moderate *unlawful* content, whereas abusive language detection research within NLP only focuses on the semantics of the content itself. Making a judgement as to whether abusive content is unlawful requires knowledge of what is legal in different countries of the world, which is not knowledge currently developed models explicitly encode. Considering just the subtask of hate speech detection can illustrate the inherent levels of nuance platforms must deal with.

**Policy challenges.** Defining a hate speech policy which is congruent with the harm that hate speech inflicts is hard. For example, Facebook's Hate Speech policy is tiered, where the Tier 1 hate speech policy aims to define hate speech that is considered most harmful. The policy covers content targeting a person or a group of persons based on their protected characteristics (race, ethnicity, national origin, disability, gender etc.). According to this policy, certain offensive statements are considered more harmful than others.

**Context.** The same content considered in different contexts can be either perfectly benign or can be hateful. For example, tagging a person in a picture of farm animals, can be perfectly benign if the person being tagged and the person doing the tagging are friends and the tagged-picture reminds them of a fond memory. On the other hand, the content could be considered hate speech if the intent is to offend the religious sensibilities of the person being tagged.

**Historical and cultural context.** Content that describes people with slurs is generally considered hate speech. However, in some cases these slur words are appropriated by the group of people who are the target of these slurs as reclaiming speech. In such cases, understanding the historical and the cultural context as well as the context of the actors involved in the speech is important.

**Content is interlinked.** People use text, images, videos, audio files, GIFs, emojis — all together — to communicate. This means that ideally, one would require a holistic understanding of the content to properly detect abusive language.

**Emerging trends and behaviors** New expressions evolve, old expressions considered acceptable may start being considered abusive in the future, speech and gestures are coopted by movements that are considered hateful, and sensibilities adapt.

The above problems are not exhaustive, but illustrate the main challenges faced by social media platforms, but new issues are regularly appearing, policies are evolving, and new detection technologies are constantly developed.

**Dataset specificity** Most of publicly available abusive language detection datasets are specific to the type of abuse, to the targeted group, or to other fine-grained distinctions, as abuse is highly context-sensitive. This diversity led to the creation of many datasets, most of them of reduced size, and much of the research on abusive language detection is fragmented. Few studies on abusive language detection view the problem from a holistic angle.

An inherent feature of hierarchical annotation schemes is that the lower levels of the taxonomy contain a subset of the instances in the higher levels, and thus there are fewer instances in the categories in each subsequent level. As more annotation levels are included in the taxonomy, there are fewer instances in each subsequent category, which makes it very difficult to train robust machine learning models.

### **5** Future Forecasting

Given the above challenges, we make the following predictions about the development of future academic research about abusive language detection:

**Focus on real-world needs.** There is a need for tighter collaboration between academia and industry to address formulations of the problem that better reflect real-world needs, e.g., focus on fine-grained categories which align with specific points in the policies of online platforms.

Multi-modality and real-time execution at scale. Given the growth of multimedia content, there will be increasing demand for moderating text within or associated with images, audio streams, and videos, which also poses challenges related to scalability and to the need for real-time inference.

**Context awareness.** In many cases, deciding whether a piece of content is abusive requires deep understanding of not just the broader context in which it is present, but a deep and broad understanding of the world, including both current and historical knowledge.

**Transparency, explainability, avoiding biases.** It is becoming increasingly crucial for a system to explain why a given piece of content was flagged: to users, to human moderators, to platform owners, to external regulatory authorities, etc. Moreover, there is a need to detect and eliminate or mitigate potential biases in the system, both due to training data or to the underlying models.

Adaptability and never-ending learning. Each platform and each community has its own policies, and different levels of tolerance regarding various types of potentially abusive content. Thus, while a universal solution could be a great start, for optimal results, the system should adapt to the target platform. As policies, moderation standards, and regulations evolve over time, there is a need for constant feedback. From a research perspective, this requires models which can easily intake new examples and feedback on the fly without the need for costly retraining.

**Federated learning.** Finally, in a cross-platform industrial setting, it is crucial to have the ability to learn from multiple platforms, while at the same time protecting user privacy.

### 6 Conclusion

We discussed recent methods for abusive language detection in the context of the content policies of online platforms. We argued that current research on abusive language detection, while a step in the right direction, does not result in the necessary tools for automating content moderation in online platforms. We argued that there is currently a dichotomy between the types of abusive language online platforms seek to curb and existing research efforts to automatically detect that abusive language. We then discussed lessons learned and major

challenges that need to be overcome. Finally, we suggested several research directions, which we forecast will emerge in the near future.

#### References

- [Adewole et al., 2017] Kayode Sakariyah Adewole et al. Malicious accounts: Dark of the social networks. *JNCA*, 79, 2017.
- [Almaatouq et al., 2016] Abdullah Almaatouq et al. If it looks like a spammer and behaves like a spammer, it must be a spammer: analysis and detection of microblogging spam accounts. *Int. J. Inf. Secur.*, 15(5), 2016.
- [Aroyehun and Gelbukh, 2018] S. T. Aroyehun and A. Gelbukh. Aggression detection in social media: Using deep neural networks, data augmentation, and pseudo labeling. In TRAC, 2018.
- [Basile et al., 2019] Valerio Basile et al. Semeval-2019 Task 5: Multilingual Detection of Hate Speech Against Immigrants and Women in Twitter. In *SemEval*, 2019.
- [Bassignana et al., 2018] Elisa Bassignana et al. Hurtlex: A multilingual lexicon of words to hurt. In *CLiC-it*, 2018.
- [Buckels *et al.*, 2014] Erin E Buckels, Paul D Trapnell, and Delroy L Paulhus. Trolls just want to have fun. *Personality and individual Differences*, 67, 2014.
- [Burnap and Williams, 2015] Pete Burnap and Matthew L Williams. Cyber hate speech on Twitter: An application of machine classification and statistical modeling for policy and decision making. *Policy & Internet*, 7(2), 2015.
- [Cambria et al., 2010] Erik Cambria et al. Do not feel the trolls. In *SDoW*, 2010.
- [Çöltekin, 2020] Çağrı Çöltekin. A corpus of Turkish offensive language on social media. In *LREC*, 2020.
- [Chen et al., 2012] Ying Chen et al. Detecting offensive language in social media to protect adolescent online safety. In PAS-SAT/SocialCom, 2012.
- [Cho et al., 2014] Kyunghyun Cho et al. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *EMNLP*, 2014.
- [Chung et al., 2019] Y. Chung et al. CONAN COunter NArratives through nichesourcing: a multilingual dataset of responses to fight online hate speech. In *ACL*, 2019.
- [Conneau et al, 2020] Conneau et al. Unsupervised cross-lingual representation learning at scale. In *ACL*, 2020.
- [Cortes and Vapnik, 1995] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3), 1995.
- [Cresci et al., 2015] Stefano Cresci et al. Fame for sale: efficient detection of fake twitter followers. *Decision Support Systems*, 80, 2015.
- [Dadvar et al., 2013] Maral Dadvar et al. Improving cyberbullying detection with user context. In *ECIR*, 2013.
- [Dave et al., 2003] Kushal Dave et al. Mining the peanut gallery: Opinion extraction and semantic classification of product reviews. In WWW, 2003.
- [Davidson et al., 2017] Thomas Davidson et al. Automated hate speech detection and the problem of offensive language. In *ICWSM*, 2017.
- [Devlin et al., 2019] Jacob Devlin et al. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.

- [Dinakar et al., 2011] Karthik Dinakar et al. Modeling the detection of textual cyberbullying. In *The Social Mobile Web*, 2011.
- [Djuric et al., 2015] Nemanja Djuric et al. Hate speech detection with comment embeddings. In *WWW*, 2015.
- [Fire et al., 2014] Michael Fire et al. Friend or foe? fake profile identification in online social networks. *Social Network Analysis and Mining*, 4(1), 2014.
- [Fortuna and Nunes, 2018] Paula Fortuna and Sérgio Nunes. A Survey on Automatic Detection of Hate Speech in Text. *ACM Computing Surveys (CSUR)*, 51(4), 2018.
- [Fortuna et al., 2018] Paula Fortuna et al. Merging datasets for aggressive text identification. In *TRAC*, 2018.
- [Founta et al., 2018] Antigoni Founta et al. Large scale crowd-sourcing and characterization of Twitter abusive behavior. In *AAAI*, 2018.
- [Founta et al., 2019] A. Founta et al. A unified deep learning architecture for abuse detection. In *WebSci*, 2019.
- [Fukushima, 1980] Kunihiko Fukushima. Neocognitron: A selforganizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4), 1980.
- [Galán-García et al., 2014] Patxi Galán-García et al. Supervised machine learning for the detection of troll profiles in twitter social network: Application to a real case of cyberbullying. In *SOCO-CISIS-ICEUTE*, 2014.
- [Gan et al., 2020] Zhe Gan et al. Large-scale adversarial training for vision-and-language representation learning. In *NeurIPS*, 2020.
- [Georgakopoulos et al., 2018] Spiros V Georgakopoulos et al. Convolutional neural networks for toxic comment classification. arXiv preprint arXiv:1802.09957, 2018.
- [Glavaš et. al, 2020] Goran Glavaš et. al. XHate-999: Analyzing and detecting abusive language across domains and languages. In *COLING*, 2020.
- [Herring et al., 2002] Susan Herring et al. Searching for safety online: Managing "trolling" in a feminist forum. *The Information Society*, 18(5), 2002.
- [Hu and Liu, 2004] Minqing Hu and Bing Liu. Mining and summarizing customer reviews. In *KDD*, 2004.
- [Juuti et al, 2020] Mika Juuti et al. A little goes a long way: Improving toxic language classification despite data scarcity. In *EMNLP Findings*, 2020.
- [Kiela et al., 2020] Douwe Kiela et al. The hateful memes challenge: Detecting hate speech in multimodal memes. *arxiv*:2005.04790, 2020.
- [Kumar et al., 2014] Srijan Kumar et al. Accurately detecting trolls in slashdot zoo via decluttering. In ASONAM, 2014.
- [Kumar et al., 2018] Ritesh Kumar et al. Benchmarking aggression identification in social media. In *TRAC*, 2018.
- [Kumar et al., 2020] Ritesh Kumar et al. Evaluating Aggression Identification in Social Media. In *TRAC*, 2020.
- [Kwok and Wang, 2013] Irene Kwok and Yuzhou Wang. Locate the hate: Detecting tweets against blacks. In AAAI, 2013.
- [Lan et al., 2019] Zhenzhong Lan et al. ALBERT: A lite BERT for self-supervised learning of language representations. *arXiv* preprint arXiv:1909.11942, 2019.

- [Li et al., 2006] Wenbin Li et al. Combining multiple email filters based on multivariate statistical analysis. In *Foundations of Intelligent Systems*, 2006.
- [Li et al., 2019] Liunian Harold Li et al. Visualbert: A simple and performant baseline for vision and language. *arxiv:1908.03557*, 2019.
- [Li et al., 2020] Xiujun Li et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In ECCV, volume 12375, 2020.
- [Lippe et al., 2020] Phillip Lippe et al. A multimodal framework for the detection of hateful memes. *arxiv*:2012.12871, 2020.
- [Majumder et al., 2018] Prasenjit Majumder et al. Filtering aggression from the multilingual social media feed. In *TRAC*, 2018.
- [Malmasi and Zampieri, 2017] Shervin Malmasi and Marcos Zampieri. Detecting Hate Speech in Social Media. In RANLP, 2017.
- [Malmasi and Zampieri, 2018] Shervin Malmasi and Marcos Zampieri. Challenges in Discriminating Profanity from Hate Speech. *Journal of Experimental & Theoretical Artificial Intelligence*, 30, 2018.
- [Mandl et al., 2019] Thomas Mandl et al. Overview of the Hasoc Track at Fire 2019: Hate Speech and Offensive Content Identification in Indo-european :anguages. In *FIRE*, 2019.
- [Mccord and Chuah, 2011] Michael Mccord and M Chuah. Spam detection on twitter using traditional classifiers. In *International Conference on Autonomic and Trusted Computing*. Springer, 2011.
- [Mittos et al., 2020] Alexandros Mittos et al. "and we will fight for our race!" A measurement study of genetic testing conversations on Reddit and 4chan. In *ICWSM*, 2020.
- [Mubarak et al., 2020] Hamdy Mubarak et al. Arabic Offensive Language on Twitter: Analysis and Experiments. *arXiv preprint arXiv:2004.02192*, 2020.
- [Muennighoff, 2020] Niklas Muennighoff. Vilio: State-of-the-art visio-linguistic models applied to hateful memes. arxiv:2012.07788, 2020.
- [Ousidhoum et al., 2019] Nedjma Ousidhoum et al. Multilingual and multi-aspect hate speech analysis. In *EMNLP-IJCNLP*, 2019.
- [Peters et al., 2018] Matthew E Peters et al. Deep contextualized word representations. In *NAACL-HLT*, 2018.
- [Pitenis et al., 2020] Zeses Pitenis et al. Offensive language identification in Greek. In *LREC*, 2020.
- [Radford et al., 2019] A. Radford et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8), 2019.
- [Ranasinghe and Zampieri, 2020] T. Ranasinghe and M. Zampieri. Multilingual offensive language identification with cross-lingual embeddings. In *EMNLP*, 2020.
- [Rosa et al, 2019] H. Rosa et al. Automatic cyberbullying detection: A systematic review. Computers in Human Behavior, 93, 2019.
- [Rosenthal et al., 2020] Sara Rosenthal et al. A large-scale semisupervised dataset for offensive language identification. *arXiv* preprint arXiv:2004.14454, 2020.
- [Rumelhart et al., 1986] David E Rumelhart et al. Learning representations by back-propagating errors. *Nature*, 323(6088), 1986.
- [Sabat et al., 2019] Benet Oriol Sabat et al. Hate speech in pixels: Detection of offensive memes towards automatic moderation. *CoRR*, abs/1910.02334, 2019.

- [Salawu et al., 2020] S. Salawu et al. Approaches to automated detection of cyberbullying: A survey. *TAFFC*, 11(1), 2020.
- [Sarna and Bhatia, 2017] Geetika Sarna and MPS Bhatia. Content based approach to find the credibility of user in social networks: an application of cyberbullying. *IJMLC*, 8(2), 2017.
- [Schmidt and Wiegand, 2017] Anna Schmidt and Michael Wiegand. A survey on hate speech detection using natural language processing. In *SocialNLP*, 2017.
- [Seah et al., 2015] Chun-Wei Seah et al. Troll detection by domain-adapting sentiment analysis. In *FUSION*, 2015.
- [Sebastiani, 2002] Fabrizio Sebastiani. Machine learning in automated text categorization. *CSUR*, 34(1), 2002.
- [Sezer et al., 2015] Baris Sezer et al. Cyber bullying and teachers' awareness. *Internet Research*, 25(4), 2015.
- [Sigurbergsson and Derczynski, 2020] Gudbjartur Ingi Sigurbergsson and Leon Derczynski. Offensive language and hate speech detection for Danish. In *LREC*, 2020.
- [Su et al., 2020] Weijie Su et al. VL-BERT: pre-training of generic visual-linguistic representations. In *ICLR*, 2020.
- [Tan and Bansal, 2019] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *EMNLP-IJCNLP*, 2019.
- [Vaswani et al., 2017] Ashish Vaswani et al. Attention is all you need. In *NIPS*, volume 30, 2017.
- [Vidgen and Derczynski, 2021] Bertie Vidgen and Leon Derczynski. Directions in abusive language training data, a systematic review: Garbage in, garbage out. *PLOS ONE*, 15(12), 2021.
- [Waseem et al., 2017] Zeerak Waseem et al. Understanding Abuse: A Typology of Abusive Language Detection Subtasks. In ALW, 2017.
- [Wiegand et al., 2018] Michael Wiegand et al. Overview of the GermEval 2018 shared task on the identification of offensive language. In *GermEval*, 2018.
- [Wong et al., 2018] Randy Y. M. Wong et al. Does gender matter in cyberbullying perpetration? An empirical investigation. *Computers in Human Behavior*, 79, 2018.
- [Wulczyn et al., 2017] Ellery Wulczyn et al. Ex machina: Personal attacks seen at scale. In *WWW*, 2017.
- [Xu et al., 2012] Jun-Ming Xu et al. Learning from bullying traces in social media. In *NAACL-HLT*, 2012.
- [Yu et al, 2020] Fei Yu et al. ERNIE-ViL: Knowledge enhanced vision-language representations through scene graph. *arxiv:2006.16934*, 2020.
- [Zampieri et al., 2019a] Marcos Zampieri et al. Predicting the type and target of offensive posts in social media. In NAACL-HLT, 2019.
- [Zampieri et al., 2019b] Marcos Zampieri et al. SemEval-2019 task 6: Identifying and categorizing offensive language in social media (OffensEval). In SemEval, 2019.
- [Zampieri et al., 2020] Marcos Zampieri et al. SemEval-2020 task 12: Multilingual offensive language identification in social media (OffensEval 2020). In *SemEval*, 2020.
- [Zhang et al., 2020] Weibo Zhang et al. Hateful memes detection via complementary visual and linguistic networks. arxiv:2012.04977, 2020.
- [Zhou et al., 2019] Luowei Zhou et al. Unified vision-language pretraining for image captioning and VQA. *arxiv:1909.11059*, 2019.

[Zhu, 2020] Ron Zhu. Enhance multimodal transformer with external label and in-domain pretrain: Hateful meme challenge winning solution. *arxiv:2012.08290*, 2020.