# UC Santa Barbara
## UC Santa Barbara Electronic Theses and Dissertations

**Title**

Text Detoxification in Natural Language Processing

**Permalink**

**Author**

Qian, Jing

**Publication Date**

2022

Peer reviewed|Thesis/dissertation

University of California
Santa Barbara

# Text Detoxification in Natural Language Processing

A dissertation submitted in partial satisfaction
of the requirements for the degree

Doctor of Philosophy
in
Computer Science

by

Jing Qian

Committee in charge:

Professor Xifeng Yan, Chair
Professor Elizabeth Belding
Professor Amr El Abbadi
Professor Heng Ji

September 2022

The Dissertation of Jing Qian is approved.

_____

Professor Elizabeth Belding

_____

Professor Amr El Abbadi

_____

Professor Heng Ji

_____

Professor Xifeng Yan, Committee Chair

August 2022

Text Detoxification in Natural Language Processing

Copyright © 2022

by

Jing Qian

I dedicate this dissertation to my husband, Peiyang, who has been my constant support through all the challenges in my life.

# Acknowledgements

I would like to start by thanking my family for their endless support throughout this journey. To my husband, thank you for sharing all my emotions, believing in me no matter what happened, and encouraging me in times of need. To my parents, thank you for your unconditional love.

I can not thank enough my advisor, Xifeng Yan. He not only gives me valuable advice in my research and career but also provides emotional support during difficult times. He is not only a role model for good researchers but also a role model for good mentors.

To Amr Elabbadi and Heng Ji, thank you for your constructive comments, which inspired me not only in this work but also in future directions. I enjoyed every discussion with you.

I would like to especially thank Elizabeth Belding for her insightful feedback and timely support. We collaborated on multiple publications and without her help, I would not have been able to quickly acquire the skills needed for academic research.

To William Wang, thank you for bringing me to UCSB and introducing me to Natural Language Processing and Computational Social Science. He gave me the foundation to start with in the early stages of this journey.

To my collaborators: Mai ElSherief, Hong Wang, Shiyang Li, Zekun Li, Ray Oshikawa, Anna Bethke, Yinyin Liu, Li Dong, Yelong Shen, Furu Wei, Weizhu Chen, Yibin Liu, Lemao Liu, Shuming Shi. Thank you for contributing to this dissertation in many ways and I have learned a lot from your wealth of knowledge.

Finally, I would like to thank all my friends for their company and encouragement. I enjoyed life in Santa Barbara because of them. They have been always supportive and understanding to me when I faced various odds.

# Curriculum Vitæ
Jing Qian

## Education

| | |
|---|---|
| 2017-2022 | Ph.D. in Computer Science, University of California, Santa Barbara. |
| | Advisor: Xifeng Yan |
| | Thesis: Text Detoxification in Natural Language Processing |
| 2013-2017 | B.Sc. in Computer Science & Technology, Nanjing University, China. |
| | Advisor: Yuzhong Qu |

## Research Interests

Computational Social Science, Deep Learning, Natural Language Generation, Natural Language Processing

## Awards & Honors

| | |
|---|---|
| 2017 | The Chancellor's Fellowship (University of California, Santa Barbara) |
| 2015 & 2016 | National Elite Program Scholarship (Nanjing University) |
| 2014 | National Scholarship (Ministry of National Education, China) |

## Publications

[1] J. Qian, L. Dong, Y. Shen, F. Wei, W. Chen. Controllable Generation with Frozen Language Models. in *Findings of the Association for Computational Linguistics (ACL 2022).*

[2] J. Qian, Y. Liu, L. Liu, Y. Li, H. Jiang, H. Zhang, S. Shi. Fine-grained Entity Typing without Knowledge Base. in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP 2021).*

[3] J. Qian, H. Wang, M. ElSherief, X. Yan. Lifelong Learning of Hate Speech Classification on Social Media. in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL 2021).*

[4] A. Gaut, T. Sun, S. Tang, Y. Huang, J. Qian, M. ElSherief, J. Zhao, D. Mirza, E. Belding, K. Chang, W. Yang Wang. Towards Understanding Gender Bias in Relation Extraction. in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL 2020).*

[5] R. Oshikawa, J. Qian, and W. Yang Wang. A Survey on Natural Language Processing for Fake News Detection. in *Proceedings of 12th International Conference on Language Resources and Evaluation (LREC 2020).*

[6] J. Qian, A. Bethke, Y. Liu, E. Belding and W. Yang Wang. A Benchmark Dataset for Learning to Intervene in Online Hate Speech. in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP 2019).*

[7] J. Qian, M. ElSherief, E. Belding and W. Yang Wang. Learning to Decipher Hate Symbols. in *Proceedings of the 2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019).*

[8] J. Qian, M. ElSherief, E. Belding and W. Yang Wang. Hierarchical CVAE for Fine-Grained Hate Speech Classification. in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP 2018).*

[9] J. Qian, M. ElSherief, E. Belding, and W. Yang Wang. Leveraging Intra-User and Inter-User Representation Learning for Automated Hate Speech Detection. in *Proceedings of The 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT 2018).*

**Experience**

**NLC Group at Microsoft Research Lab - Asia**, 10/2020-07/2021
Research Intern, Advisor: Li Dong, Yelong Shen, Furu Wei, and Weizhu Chen
**NLP Group at Tencent AI Lab**, 06/2020-09/2020
Research Intern, Advisor: Lemao Liu
**Deep Learning Group, Microsoft Research AI**, 06/2019-09/2019
Research Intern, Advisor: Sungjin Lee
**CS 190I/291A, Special Topics in Computer Science**, UCSB, Fall 2018
Teaching Assistant
**University of Waterloo**, 09/2016-04/2017
Non-degree program, Computer Science
**NJU Websoft Group, Nanjing University**, 06/2015-07/2017
Undergraduate Research Assistant, Advisor: Yuzhong Qu

# Abstract

Text Detoxification in Natural Language Processing

by

Jing Qian

The rapid rise of social media has given individuals an easy platform to publicly communicate with others. Unfortunately, this has also led to improper use of online spaces, such as the propagation of toxic speech. Worse yet, user-generated toxic speech can propagate beyond online social platforms. We recently observed toxic degeneration and biased behavior in language models (LMs) pretrained on a web text corpus. In order to create a healthy communication environment, developing text detoxification techniques has attracted increasing attention from both the industry and academia. To achieve this goal, detoxification methods need to handle both pre-existing toxic speech and LMs that exhibit toxic degeneration.

In this dissertation, we investigate how to use Natural Language Processing techniques for text detoxification in two directions. 1) for pre-existing toxic speech, we develop automatic tools for post-processing, including detection, analysis, and intervention. 2) for the model-generated toxic speech, another complementary solution to post-processing is to detoxify the pretrained LMs by reducing the likelihood that the model will generate toxic content.

In the first part, we focus on toxic speech, especially hate speech, that already exists. We start by improving automated hate speech detection through intra-user and inter-user representation learning. We then move beyond the standard binary hate speech detection. We study fine-grained hate speech characterizing in both the isolated learning setting and the lifelong learning setting. We also investigate how neural network models

are able to decipher hate symbols. We then explore the intervention strategies for online conversations that contain hate speech. As part of this work, we make publicly available two fully-labeled hate speech datasets with human-written intervention responses.

In the next part, we focus on the not-yet-emergent toxic speech from LMs. We begin by controlling pretrained LMs in the freeform text generation scenario. We then further investigate LM detoxification in dialogue with contextualized stance control. Our work effectively lowers the toxic content rate of the pretrained LMs while sacrificing less linguistic quality.

Finally, we summarize the key findings of our work and discuss future research directions to push the boundaries of text detoxification.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

The rapid growth of social media has not only yielded a vast increase in information accessibility but has also given individuals an easy platform on which to share their beliefs and publicly communicate with others. Unfortunately, the growing popularity of online interactions through social media has also led to nefarious uses of online spaces, such as the propagation of toxic speech, which is *impolite, rude or hurtful language.* In toxic speech, hate speech specifically refers to *language that attacks others on the basis of the protected characteristics, such as race and religion* [1]. For a healthier communication environment, we aim to detoxify text using Natural Language Processing techniques, as shown in Figure 1.1.

Developing automated text detoxification techniques has realistic significance. Toxic speech can cause psychological harm to its victims. Research in psychology finds that being targeted by online hate speech leads to greater stress expression [2]. To reduce the propagation of toxic speech, social media platforms such as Facebook hire content moderators to review content, but it is shown that repeated exposure to these traumatizing posts has left some content moderators with mental health problems [3]. Therefore, automatic text detoxification can not only reduce the exposure of the target people to

Figure 1.1: Text Detoxification.

toxic speech but also reduce the content moderators' exposure to it. Besides, hate speech analysis may provide valuable information about real-life hate crimes. A recent study has found a connection between the number of real-life racially motivated hate crimes and the number of racist messages posted on Twitter in 100 different cities across the United States [4].

In the first part of the dissertation, we apply post-processing to tackle toxic speech, especially hate speech, that already exists. We begin with automatic hate speech detection, which differentiates hate speech content from normal content. Most previous work on hate speech analysis using Natural Language Processing (NLP) techniques falls into this category [5, 6, 1]. We then investigate fine-grained hate speech characterization, which is necessary for fine-grained hate speech analysis. Different from hate speech detection, the input data for hate speech characterization constitute hate speech only and the output can be the target of hate speech, the likely spread of hate speech, etc. In this work, we are particularly interested in analyzing hate groups and hate symbols. Towards

Figure 1.2: Dissertation Overview.

the end goal of mitigating the problem of online toxic speech, we then explore methods to automatically generate responses to intervene during online conversations containing hate speech.

In the second part, we tackle the toxic degeneration problem of Language Models. Language models pretrained on problematic web content are significantly at risk of producing toxic content [7], which hinders the safe deployment of these powerful Natural Language Generation (NLG) models. We propose novel control methods to alleviate this problem in the freeform generation scenario as well as in a more complex scenario involving dialogue contexts. Figure 1.2 shows an overview of the dissertation.

To conclude, we demonstrate that:

*Countering toxic speech is a critical yet challenging task, but one which can be aided by the use of Natural Language Processing techniques. We advance the state of the art in automatic post-processing of hate speech and Language Model detoxification by leveraging both Natural Language Understanding and Natural Language Generation techniques.*

| Source | Definition |
|---|---|
| European Union Commission [8] | *All conduct publicly inciting to violence or hatred directed against a group of persons or a member of such a group defined by reference to race, colour, religion, descent or national or ethnic origin.* |
| Nobata et al. [1] | *Language which attacks or demeans a group based on race, ethnic origin, religion, disability, gender, age, or sexual orientation/gender identity.* |
| Twitter [9] | *You may not promote violence against or directly attack or threaten other people on the basis of race, ethnicity, national origin, sexual orientation, gender, gender identity, religious affiliation, age, disability, or disease.* |

Table 1.1: Definitions of hate speech from different sources. The targets of hate speech, the characteristics that hate speech is based on, and the intent of hate speech are marked in different colors.

## 1.1    Hate Speech and Toxic Speech

As shown in Table 1.1, there are some nuances between definitions of hate speech from different sources, but they converge on three dimensions: specific targets, specific characteristics, and the intent to attack or incite violence. First, hate speech has specific targets. It can be an individual or a group of people. Second, hate speech is based on the protected characteristics of the target, such as race, ethnic origin, etc. Third, the intent of hate speech is to attack or promote violence. Figure 1.3 shows two examples of hate speech on social media. The first example targets people with a specific disease while the second example targets a religious group.

Compare to hate speech, toxic speech is a more general term referring to impolite, rude, or hurtful language. There does not exist a formal definition of the distinction between hate speech and toxic speech but there is a consensus that hate speech targets disadvantaged social groups or group members in a manner that is potentially harmful to them while toxic speech may not [10, 11, 12]. In Computer Science research [10,

Figure 1.3: Examples of hate speech based on disease (left) and hate speech based on religion (right). The names of the disease and the religion are masked.

7, 1], the terms, toxic speech, offensive language, and abusive language, are often used interchangeably. In the first part of the dissertation, we specifically focus on the issue of hate speech. In the second part of the dissertation, we focus on the issue of more general toxic speech.

## 1.2  Language Models

A language model is the core component of Natural Language Generation (NLG). Given a text consisting of $m$ tokens $(x_1, x_2, \ldots, x_m)$, a language model assigns a probability $P(x_{1:m}) = P(x_1, x_2, \ldots, x_m)$ to the text. Typically, the joint probability $P(x_{1:m})$ is decomposed as $\prod_{i=1}^{m} P(x_i|x_{1:i-1})$ and a neural network model is used to estimate the conditional probability $P(x_i|x_{1:i-1})$. Such neural network models are called neural language models. As illustrated in Figure 1.4, the neural language model estimates the probability distribution of the next word token given the text "The book is". The next token "well" is selected or sampled based on this estimated distribution, so now we have the text "This book is well". The above steps are repeated until we reach the end of the sentence. Therefore, with a language model, natural language generation can be carried out naturally, and the quality of the language model is critical to the quality of the generated natural language text.

Recently, the advent of Transformer [13] enables large LMs and large-scale pretraining

Figure 1.4: An example of Transformer-based Natural Language Generation Models.

[14, 15, 16]. Figure 1.4 shows the structure of GPT2 [14], a large Transformer-based Language Model. It stacks multiple layers of the Transformer decoder, which consists of a masked self-attention followed by a feed-forward layer. During pretraining, the LM learns to estimate the conditional probability $P(x_i|x_{1:i-1})$ from a large unlabeled corpus. Currently, the prevalent LMs are pretrained on large amounts of text from the web, such as news, blogs, and so on. Large-scale pretraining enables these LMs to achieve a good estimation of the next token probability, resulting in coherent, human-like texts. However, relying on web text for pretraining also results in degeneration and biased behavior. Previous research [7] shows that the language models pretrained on problematic web content are significantly at risk of producing toxic content. These pretrained LMs can easily degenerate into toxicity, even without explicitly toxic prompts. This problem hinders the safe deployment of these powerful LMs for Natural Language Generation.

## 1.3    Text Detoxification

To tackle toxic speech from both online users and pretrained language models, one solution is to do post-processing. A classifier can be trained to automatically filter out the toxic content. An extensive body of work has focused on this field. Most of the previous work focuses on binary classification [5, 6, 1, 17] or coarse-grained multi-class classification [18, 19, 10]. In the case of user-generated toxic content, the detected content can be removed and further analysis or intervention can be performed [20, 21, 22, 23]. In the case of model-generated toxic content, in addition to discarding the toxic content, the NLG model needs to repeatedly do the generation until the classifier finds a safe one. However, this strategy is inefficient and there exist cases where no safe choice exists within a fixed number of generations. Previous research [24] finds a universal trigger of hate speech for GPT2. The authors use this trigger to generate 100 samples from GPT2 and manually analyze them. They find 98% are racial or offensive. In this case, post-processing may not work effectively.

In order to circumvent this limitation, some research [25, 26, 27] has focused on controlling the generation of these Transformer-based LMs from the source, where the desired attributes are directly infused into the generation process. Detoxification can be considered as an application of the control method to infuse the safe language attribute. By controlling a generation model, we would like to achieve two goals. On one hand, we would like the generations to have good attribute alignment. If the desired attribute is non-toxic, the control method should guide the model to avoid toxic generations. On the other hand, we would like to maintain good linguistic quality. As mentioned before, a large pretrained language model can produce coherent texts. When it is augmented with a control method, we do not want to sacrifice the linguistic quality too much. If the control method causes the generation to be cluttered, the advantage of using the large

pretrained language model is lost. Therefore, good attribute alignment and linguistic quality should be achieved at the same time.

## 1.4    Contributions

In the first part of the dissertation, we focus on automatic post-processing of hate speech. For **hate speech detection**, we provide a novel perspective on hate speech detection by integrating target tweets, intra-user, and inter-user representations in a unified framework, which outperforms strong baselines [17]. For **hate speech analysis**, we introduce the first work on fine-grained hate speech characterization that attributes hate groups to each tweet, and our proposed Hierarchical Conditional VAE model improves the Micro-F1 score by up to 10% over the baselines [20]. Based on this work, we further introduce the first work on lifelong learning of fine-grained hate speech characterization and propose a novel method that utilizes Variational Representation Learning along with a memory module to alleviate catastrophic forgetting. Experimental results show that our proposed method outperforms the state-of-the-art significantly on the average F1 scores [21]. In addition to the standard formulation of hate speech classification, we also propose a novel task of learning to decipher hate symbols. We investigate sequence-to-sequence neural network models along with a novel Variational Decipher and show how they are able to crack the hate symbols [22]. For **hate speech intervention**, we introduce the generative hate speech intervention task and provide two fully-labeled hate speech datasets with human-written intervention responses. The released datasets and our study on them provide a benchmark for future research in this field [23].

In the second part of the dissertation, we focus on language model detoxification. For **LM detoxification in freeform text generation**, we propose a novel lightweight control method that utilizes contrastive prefixes along with frozen LMs. We propose a

supervised method and an unsupervised method with novel objectives for prefix training and experimental results show that our methods can effectively detoxify the LMs while keeping high linguistic quality [27]. Based on this work, for **LM detoxification in dialogue response generation**, we further propose a novel control framework that combines context-dependent and context-independent control by utilizing hierarchical prefixes and novel contrastive training objectives. Experimental results show that our proposed method can effectively learn the contextualized stance control while keeping a low self-toxicity of the generation model.

With these efforts, we aim to reduce the spread of toxic speech and its harm to others by developing automated text detoxification methods.

# Part I

# Automatic Post-processing of Hate Speech

# Chapter 2

# Hate Speech Detection

In this chapter, we improve automated hate speech detection by presenting a novel model that leverages intra-user and inter-user representation learning for robust hate speech detection on Twitter. In addition to the target tweet, we collect and analyze the user's historical posts to model intra-user tweet representations. To suppress the noise in a single tweet, we also model the similar tweets posted by all other users with reinforced inter-user representation learning techniques. Experimentally, we show that leveraging these two representations can significantly improve the F-score of a strong bidirectional LSTM baseline model by 10.1%.

## 2.1 Introduction

Despite the existence of numerous studies dedicated to the development of NLP hate speech detection approaches, the accuracy is still poor. The central problem is that social media posts are short and noisy, and most existing hate speech detection solutions take each post as an isolated input instance, which is likely to yield high false positive and negative rates.

Figure 2.1: Our hate speech classifier. In contrast to existing methods that focus on a single target tweet as input (center), we incorporate intra-user (right) and inter-user (left) representations to enhance performance.

In this paper, we focus on augmenting hate speech classification models by first performing representation learning to model user history without supervision. The hypothesis is that, by analyzing a corpus of the user's past tweets, our system will better understand the language and behavior of the user, leading to better hate speech detection accuracy. Another issue is that using a single tweet as input is often noisy for any machine learning classifier. For example, the tweet *"I'm not sexist but I can not stand women commentators"* is actually an instance of hate speech, even though the first half is misleading. To minimize noise, we also consider semantically similar tweets posted by other users. To do so, we propose a reinforced bidirectional long short-term memory network (LSTM) [28] to interactively leverage the similar tweets from a large Twitter dataset to enhance the performance of the hate speech classifier. An overview of our approach is shown in Figure 2.1. The main contributions of our work are:

- We provide a novel perspective on hate speech detection by modeling intra-user tweet representations.

- To improve robustness, we leverage similar tweets from a large unlabeled corpus with reinforced inter-user representations.

- We integrate target tweets, intra-user, and inter-user representations in a unified framework, outperforming strong baselines.

## 2.2   Related Work

Most previous work in automated hate speech classification focuses on binary classification, such as differentiating between anti-semitic or not [5] and differentiating between abusive or not [6, 1]. In [29], the authors build a set of binary classifiers to detect hate speech based on race, sexual orientation, or disability. Some research focuses on three-way classification [18, 10]. In [18], the authors classify each input tweet as racist hate speech, sexist hate speech, or neither, while in [10], a model is trained to differentiate among three classes: containing hate speech, only offensive language, or neither. In this chapter, we formulate the task as binary classification, differentiating between hate speech or not.

A recent survey outlined eight categories of features used in hate speech detection [30]: simple surface [5, 18], word generalization [5, 6], sentiment analysis [31], lexical resources and linguistic features [29], knowledge-based features [32], meta-information [18], and multi-modal information [6]. Closely related to our work is research that leverages user attributes in the classification process such as the history of participation in hate speech and usage of profanity [33, 34]. Both [33] and [34] collect user history to enhance detection accuracy. In [33], the user history is required to be labeled instances. However, labeling user history requires significant human effort. In [34], the authors model the user with manually selected features. In contrast, our approach leverages unlabeled user history to automatically model the user.

Figure 2.2: The overview of our proposed model. $t$ is the input target tweet, $z$ denotes intra-user tweets, and $x_a$ is the selected inter-user tweet. $r_{ie}$ is the inter-user representation, $r_{ia}$ is the intra-user representation, and $r_{ta}$ is the representation of the target tweet. These three branches respectively correspond to the three branches illustrated in Figure 2.1. $y_i$ is the prediction at the time step $i$ and $s_i$ is the state input for the agent at the time step $i$.

## 2.3    Methods

Figure 2.2 illustrates the architecture of our model. It includes three branches, whose details will be described in the following subsections.

### 2.3.1    Bidirectional LSTM

Given a target tweet, the baseline approach is to feed the embeddings of the tweet into a bidirectional LSTM network [28, 35, 36] to obtain the prediction. This is shown in the middle branch in Figure 2.1. However, this method is likely to fail when the target tweet is noisy or the critical words for making predictions are out of vocabulary.

## 2.3.2   Intra-User Representation

The baseline approach does not fully utilize available information, such as the user's historical tweets. In our approach, we collect the user's historical posts through the Twitter API. For a target tweet $t$, suppose we collect $m$ tweets posted by this user: $Z_t = \{z_1, z_2, ..., z_m\}$. These intra-user tweets are fed into a pre-trained model to obtain an intra-user representation. The pre-trained model has the same structure as the baseline model. This is shown in the right branch in Figures 2.1 and 2.2. The intra-user representation is then combined with the baseline branch for the final prediction. The computation process is:

$$o_{ta}(t) = f_{ta}(t, \mathbf{0}) \tag{2.1}$$

$$r_{ta}(t) = l_{ta}(\sigma(o_{ta}(t))) \tag{2.2}$$

$$o_{ia}(z_j) = f_{ia}(z_j, \mathbf{0}) \tag{2.3}$$

$$r_{ia}(t) = \sigma(\sum_{j=1}^{m} l_{ia}(\sigma(o_{ia}(z_j)))) \tag{2.4}$$

where $f_{ta}$ is the bi-LSTM of the baseline branch; $o_{ta}$ is the output of the bi-LSTM; and $l_{ta}$, $l_{ia}$ are linear functions. Similarly, $f_{ia}$ is the bi-LSTM of the intra-user branch and $o_{ia}$ is the output. $r_{ta}$ is the output prediction of the baseline branch. $r_{ia}$ is the intra-user representation, and $\sigma$ is the non-linear activation function.

## 2.3.3   Inter-User Representation

In addition to the user history, the tweets that are semantically similar to the target tweet can also be utilized to suppress noise in the target tweet. We collect similar tweets from large unlabeled tweet set $U$ by Locality Sensitive Hashing (LSH) [37, 38]. Since the space of all tweets is enormous, we use LSH to efficiently reduce the search space. For each

target tweet $t$, we use LSH to collect $n$ nearest neighbors of $t$ in $U$: $x_1$, $x_2$, ..., $x_n$. These $n$ tweets form the inter-user tweet set for $t$: $X_t = \{x_1, x_2, ..., x_n\}$. Due to the size of this set, a policy gradient-based deep reinforcement learning agent is trained to interactively fetch inter-user tweets from $X_t$. The policy network consists of two layers as shown in the middle part of Figure 2.2 and the policy network is trained by the REINFORCE algorithm [39]. At each time step $i$, the action of the agent is to select one tweet $x_a$ from $X_t$. $x_a$ is then fed into a bi-LSTM followed by a linear layer. The result is combined with the intra-user representation and the baseline prediction (the right and the middle branch in Figures 2.1 and 2.2) to get the prediction at time step $i$. At each time step, the bi-LSTM layer that encodes the selected inter-user is initialized with the output hidden state of the last time step. The number of time steps for each target tweet is set to be a fixed number $T$ so that the agent will terminate after $T$ fetches. The final prediction occurs at the last time step. The computation is shown by the following equations.

$$o_{ie}(x_a)_i, h_{ie}(x_a)_i = f_{ie}(x_a, h_{ie}(x_b)_{i-1}) \tag{2.5}$$

$$r_{ie}(x_a)_i = l_{ie}(\sigma(o_{ie}(x_a)_i)) \tag{2.6}$$

$$y'(t) = \sigma(l_c(r_{ta}(t) \oplus r_{ia}(t))) \tag{2.7}$$

$$y(t)_i = \sigma(l_c(r_{ie}(t) \oplus r_{ta}(t) \oplus r_{ia}(t))) \tag{2.8}$$

where $x_a$ is the selected inter-user tweet at time step $i-1$. $f_{ie}$ is the bi-LSTM of the inter-user branch. $o_{ie}$ and $h_{ie}$ are the output and the hidden state. $l_c$ is a linear function. $r_{ie}$ is the inter-user representation. $y'$ is the prediction made without the inter-user branch and $y$ is the prediction made with the inter-user branch. The symbol $\oplus$ means concatenation. The subscript $i$ denotes time step $i$.

The state at each time step for the agent is the concatenation of encoded inter-user tweets, the output of the Bi-LSTM in the inter-user branch and the baseline branch,

together with the intra-user representation in the intra-user branch (the dotted arrows in Figure 2.2). Each inter-user tweet $x_j$ in $X_t$ is encoded by the bi-LSTM of the inter-user branch (the dotted arrow through the Bi-LSTM of the inter-user branch).

$$b(x_j) = f_{ie}(x_j, \mathbf{0}) \tag{2.9}$$

$$s(t)_i[j] = o_{ie}(x_a)_i \oplus b(x_j) \oplus o_{ta}(t) \oplus r_{ia}(t) \tag{2.10}$$

$b$ is the output of the bi-LSTM of the inter-user branch. In order to differentiate with $o_{ie}$ mentioned above, we use $b$. $s(t)_i[j]$ is the $j$th row of the state at time step $i$.

By using reinforcement learning, the state for the agent is updated after each fetch of the inter-user tweet. Thus, the agent can interactively make selections and update the inter-user representations step by step. The reward $v_i$ for the agent at time step $i$ is based on the original prediction without the agent and the prediction at the last time step with the agent. The computation is shown as:

$$q(t)_i = e(y'(t), y^*) - e(y(t)_T, y^*) \tag{2.11}$$

$$v(t)_i = \begin{cases} \alpha * q(t)_i & \text{if } y'(t)! = y^* \\ q(t)_i & \text{else if } y(t)_T! = y^* \\ 0 & \text{otherwise} \end{cases} \tag{2.12}$$

where $e$ is the loss function; $q(t)$ is the basic reward; and $v(t)_i$ is the modified reward at time step $i$. $\alpha$ is a positive number used to amplify the reward when the original classification is incorrect. The intuition of this reward is to train the agent to be able to correct the misclassified tweets. When the original prediction and the last prediction are both correct, the reward is set to 0 to make the agent focus on the misclassified instances.

The complete training process is shown in Algorithm 1. Before the training, the intra-

---

**Algorithm 1** Training Algorithm

---

1: **for** $t$ in training set **do**
2:       collect $X_t$ and $Z_t$;
3:       compute intra-user representation $r_{ia}(t)$;
4: **end for**
5: initialize parameters $\theta_p$ of the policy network;
6: initialize parameters $\theta_e$ of the other nets;
7: **for** $epoch = 1, E$ **do**
8:       **for** $t$ in training set **do**
9:             compute $o_{ta}(t)$, $r_{ta}(t)$;
10:             compute the raw prediction $y'(t)$;
11:             compute $b(X_t)$;
12:             $x_a = t$;
13:             compute $o_{ie}(t)$;
14:             initialize the state $s(t)_0$;
15:             **for** $i = 1, T$ **do**
16:                   agent select action by $\epsilon - greedy$;
17:                   update $x_a$;
18:                   compute $o_{ie}(t)$, $r_{ie}(t)$;
19:                   compute $y(t)_i$ and $s(t)_i$;
20:                   compute the reward $v_i(t)$;
21:             **end for**
22:             apply REINFORCE to update $\theta_p$;
23:             update $\theta_e$ on the loss $\mathcal{L}(\theta_e) = e(y(t)_T, y^*)$;
24:       **end for**
25: **end for**

---

user tweets and inter-user tweets are collected for each target tweet. Then intra-user representations are computed, followed by the computation for initializing the environment and state for the agent. Next, the agent's actions, state updates, prediction, and reward are computed. Finally, the parameters are updated.

18

## 2.4   Experiments and Discussion

### 2.4.1   Experimental Settings

**Data:** We use the English hate speech dataset published in [18]. This dataset consists of 16,907 annotated tweets. The original dataset only contains the tweet ID and the label for each tweet. We expand the dataset with the user ID and tweet text. After deleting the tweets that are no longer accessible, the dataset we use contains 15,781 tweets from 1,808 users. In the published dataset, each tweet is manually labeled as racism, sexism, and none. Since we consider a binary classification setting, we union the first two sets. In the final dataset, 67% are labeled as non-hate speech, and 33% are labeled as hate speech. 1000 tweets are randomly selected for testing and the remaining 14,781 tweets are for training.

**Measures:** We use F1 scores for evaluation. The F1 score is defined as the harmonic mean of precision and recall, so we also report the precision and recall scores achieved by each method in the experiments.

**Baseline:** The baseline model is a bi-LSTM model. The input for the model is the word embeddings of the target tweet. The word embedding is of size 200. The hidden size of the LSTM is 64. The optimizer is Adam and we use mini-batches of size 25. The word embedding model is pre-trained on a tweet corpus containing 3,433,513 tweets.

**Intra-user Representation Learning:** Based on the target tweet, we collect at most 400 tweets posted by the same user, with the target tweet removed. The baseline branch and the intra-user branch are combined via a linear layer.

**Combining with Inter-user Representation:** The inter-user tweet set is collected from the dataset via Locality Sensitive Hashing (LSH). In our experiments, we use a set size of either 50, 100, or 200 tweets. At each time step, one tweet is selected from the inter-user tweet set by the policy agent. We also experimented with a second setting, in

| Method | Prec. | Rec. | F1 |
|---|---|---|---|
| SVM | **.793** | .656 | .718 |
| Logistic Regression | .782 | .611 | .686 |
| Bi-LSTM + attention | .760 | .665 | .710 |
| CNN-static | .701 | .707 | .703 |
| CNN-non-static | .743 | .699 | .720 |
| N-gram | .729 | **.777** | .739 |
| Bi-LSTM | .672 | .737 | .703 |
| + Intra. Rep. | .772 | .749 | .760 |
| + Intra.+ Randomized Inter. Rep. | .773 | .764 | .768 |
| + Intra.+ Reinforced Inter. Rep. | .775 | .773 | **.774** |

Table 2.1: Experimental results. Prec.: precision. Rec.: recall. F1: F measure. Bi-LSTM: the baseline bidirectional LSTM model. Bi-LSTM + attention: an attentional bidirectional LSTM model. The experimental settings of the last three rows are illustrated in Section 2.4.1. + Intra. Rep.: the model consists of the target tweet branch and the intra-user branch. + Intra. + Randomized Inter. Rep. incorporates randomly selected inter-user tweets while + Intra. + Reinforced Inter. Rep. further incorporates the reinforced inter-user branch. The best results are in bold.

which we replace the agent by random selection. At each time step, an inter-user tweet is randomly selected from $X$ and fed into the inter-user branch.

## 2.4.2   Results

We compare the above settings with six classification models: Support Vector Machine (SVM) [40], Logistic Regression, attentional BI-LSTM, two CNN models [41], and an N-gram model [18]. We evaluate these models on three metrics: precision, recall, and F1 score. The results are shown in Table 2.1. We report results for $|U| = 100$ in Table 2.1, as results with sizes 50 and 200 are similar. We find that leveraging the intra-user information helps reduce false positives. The performance is further improved when integrating our model with inter-user similarity learning. Our results show that selection by the policy gradient agent is slightly better than random selection, and we hypothesize the effect would be more salient when working with a larger unlabeled dataset. The

McNemar's test shows that our model gives significantly better (at $p < 0.01$) predictions than the baseline bi-LSTM and attentional bi-LSTM.

### 2.4.3   Error Analysis

There are two types of hate speech that are misclassified. The first type contains rare words and abbreviations, e.g. *FK YOU KAT AND ANDRE! #mkr*. Such intentional misspellings or abbreviations are highly varied, making it difficult for the model to learn the correct meaning. The second type of hate speech is satire or metaphor, e.g. *Congratulations Kat. Reckon you may have the whole viewer population against you now #mkr*. Satire and metaphors are extremely difficult to recognize. In the above two cases, both the baseline branch and the inter-user branch can be unreliable.

## 2.5   Conclusion

In this work, we propose a novel method for hate speech detection. We use bi-LSTM as the baseline method. However, our framework can easily augment other baseline methods by incorporating intra-user and reinforced inter-user representations. In addition to detecting potential hate speech, our method can be applied to help detect suspicious social media accounts. Considering the relationship between online hate speech and real-life hate actions, our solution has the potential to help analyze real-life extremists and hate groups. Furthermore, intra-user and inter-user representation learning can be generalized to other text classification tasks, where either user history or a large collection of unlabeled data is available.

# Chapter 3

# Fine-grained Hate Speech Characterization

Automated hate speech detection typically focuses on binary classification. In this chapter and the next chapter, we focus on hate speech characterization, the goal of which is to analyze hate speech instead of distinguishing hate speech from non-hate speech. We propose a novel method for a fine-grained hate speech classification task, which focuses on differentiating among hate groups of different hate group categories. We first explore the Conditional Variational Autoencoder (CVAE) [42, 43] as a discriminative model and then extend it to a hierarchical architecture to utilize the additional hate category information for more accurate prediction. Experimentally, we show that our proposed model outperforms commonly-used discriminative models.

## 3.1   Introduction

In response to the growth in online hate, there has been a trend of developing automatic hate speech detection models to alleviate online harassment [5, 18]. However, a

Figure 3.1: A portion of the hate group map published by the Southern Poverty Law Center. Each marker represents a hate group. The markers with the same pattern indicate the corresponding hate groups share the same ideology. The white box shows an example of a hate group in Auburn, Alabama under the category of "White Nationalist". The name of the group is masked.

common problem with these methods is that they focus on coarse-grained classifications with only a small set of categories. To the best of our knowledge, the existing work on hate speech detection formulates the task as a classification problem with no more than seven classes. Building a model for more fine-grained multi-class classification is more challenging since it requires the model to capture finer distinctions between each class.

Moreover, fine-grained classification is necessary for fine-grained hate speech analysis. Figure 3.1 is a portion of the hate group map published by the Southern Poverty Law Center (SPLC), where a hate group is defined as "an organization that – based on its official statements or principles, the statements of its leaders, or its activities – has beliefs or practices that attack or malign an entire class of people, typically for their immutable characteristics." The SPLC divides the 954 hate groups in the United States into 17 categories according to their ideologies (e.g. Racist Skinhead, Anti-Immigrant, and others). The SPLC monitors hate groups throughout the United States by a variety of methodologies to determine the activities of groups and individuals: reviewing hate group

publications and reports by citizens, law enforcement, field sources, and the news media, and conducting their own investigations. Therefore, building automatic classification models to differentiate between social media posts from different hate groups is both challenging and of practical significance.

In this paper, we propose a fine-grained hate speech classification task that separates tweets posted by 40 hate groups of 13 different hate group categories. Although CVAE is commonly used as a generative model, we find it can achieve competitive results and tends to be more robust when the size of the training dataset decreases, compared to the commonly used discriminative neural network models. Based on this insight, we design a Hierarchical CVAE model (HCVAE) for this task, which leverages the additional hate group category (ideology) information for training. Our contributions are three-fold:

- This is the first work on fine-grained hate speech classification that attributes hate groups to individual tweets.

- We propose a novel Hierarchical CVAE model for fine-grained tweet hate speech classification.

- Our proposed model improves the Micro-F1 score by up to 10% over the baselines.

## 3.2   Background and Related Work

### 3.2.1   Text Classification

Among the previous work in automated hate speech classification, our work is most closely related to [31], which focuses on fine-grained cyberbullying classification. However, this study only focuses on seven categories of cyberbullying while our dataset consists of dozens of classes. Therefore, our classification task is much more fine-grained and challenging.

24

Our work is also related to fine-grained text classification. Convolutional Neural Networks (CNN) have been successfully applied to the text classification task. In [41], CNN is applied at the word level while in [44], CNN is applied at the character level. In [45], the authors exploit the word order of text data with CNN for accurate text categorization. The Recursive Neural Tensor Network is introduced for text classification in previous work [46]. Recurrent Neural Networks (RNN) are also commonly used for text classification [47, 48]. Previous work further combines RNN with CNN [49, 50]. The hierarchical structure is previously exploited for document classification [51, 52]. In [51], the model generates from the sentence representation to the document representation while in [52], the model generates from the word-level attention to the sentence-level attention. However, the division of the hierarchies in our HCVAE is according to semantic levels, rather than according to document compositions. We generate from the category-level representations to the group-level representations. Moreover, the most commonly used datasets by these works (Yelp reviews, Yahoo answers, AGNews, IMDb reviews [53]) have no more than 10 classes.

### 3.2.2   Variational Autoencoder

Variational Autoencoder (VAE) [54] is a generative model composed of both an encoder and a decoder, but unlike a sequence-to-sequence model that directly encodes the input into a latent space, VAE approximates the underlying probability distribution of data. VAE has achieved competitive results in many complicated generation tasks, such as handwritten digits [54, 55], faces [54, 56], and machine translation [57]. CVAE [42, 43] is an extension of the original VAE framework that incorporates conditions during generation. In addition to image generation, CVAE has also been successfully applied to several NLP tasks, such as dialog generation [58]. Although so far CVAE has always been

used as a generative model, we explore the performance of the CVAE as a discriminative model and further propose a hierarchical CVAE model, which exploits the hate group category (ideology) information for training.

## 3.3   Data, Task, and Measures

We collect the data from 40 hate group Twitter accounts of 13 different hate ideologies, e.g., white nationalist, anti-immigrant, racist skinhead, among others. The detailed themes and core values behind each hate ideology are discussed in the SPLC ideology section [59]. For each hate ideology, we collect a set of Twitter handles based on hate groups identified by the SPLC [60]. For each hate ideology, we select the top three handles in terms of the number of followers. Due to ties, there are four different groups in several categories of our dataset. The dataset consists of all the content (tweets, retweets, and replies) posted with each account from the group's inception date, as early as 07-2009, until 10-2017. Each instance in the dataset is a tuple of (tweet text, hate category label, hate group label). The complete dataset consists of approximately 3.5 million tweets. Note that due to the sensitive nature of the data, we anonymously reference the Twitter handles for each hate group by using IDs throughout this paper. The distribution of the data is illustrated in Figure 3.2.

As shown in Figure 3.2, the dataset is highly imbalanced, which causes problems for training the above models. In order to alleviate this problem, we use a weighted random sampler to sample instances from the training data. However, the testing dataset is not sampled, so the distribution of the testing dataset remains the same as that of the original dataset. This allows us to evaluate the models' performance on the data with a realistic distribution.

Given a tweet text, our goal is to predict which hate group out of the 40 it is from.

Figure 3.2: The distribution of the data for 40 hate groups. The X-Axis is each hate group. The Y-Axis is the number of tweets collected for each hate group.

The ground truth hate ideology of a tweet text is unavailable during testing, but can be utilized during training. We use Macro-F1, Micro-F1, and Weighted-F1 to evaluate the classification results.

## 3.4 CVAE Baseline

We formulate our classification task as the following equation:

$$Obj = \sum_{(x,y_g) \in X} \log p(y_g|x) \tag{3.1}$$

where $x$, $y_g$ are the tweet text and hate group label respectively, $X$ is the dataset. Instead of directly parameterizing $p(y_g|x)$, it can be further written as the following equation:

$$p(y_g|x) = \int_z p(y_g|z, x) p(z|x) dz \tag{3.2}$$

27

where $z$ is the latent variable. Since the integration over $z$ is intractable, we instead try to maximize the corresponding evidence lower bound (ELBO):

$$ELBO = E[\log p(y_g|z, x)] - $$
$$D_{KL}[q(z|x, y_g)||p(z|x)] \tag{3.3}$$

where $D_{KL}$ is the Kullback–Leibler (KL) divergence. $p(y_g|z, x)$ is the likelihood distribution, $q(z|x, y_g)$ is the posterior distribution, and $p(z|x)$ is the prior distribution. These three distributions are parameterized by $p_\varphi(y_g|z, x)$, $q_\alpha(z|x, y_g)$, and $p_\beta(z|x)$. Therefore, the training loss function is:

$$\mathcal{L} = L_{REC} + L_{KL}$$
$$= E_{z \sim p_\alpha(z|x, y_g)}[-\log p_\varphi(y_g|z, x)] + \tag{3.4}$$
$$D_{KL}[q_\alpha(z|x, y_g)||p_\beta(z|x)]$$

The above loss function consists of two parts. The first part $L_{REC}$ is the reconstruction loss. Optimizing $L_{REC}$ can push the predictions made by the posterior network and the likelihood network closer to the ground truth labels. The second part $L_{KL}$ is the KL divergence loss. Optimizing it can push the output distribution of the prior network and that of the posterior network closer to each other, such that during testing, when the ground truth label $y_g$ is no longer available, the prior network can still output a reasonable probability distribution over $z$.

Figure 3.3 illustrates the structure of the CVAE model during training (the structure used for testing is different). The likelihood network $p_\varphi(y_g|z, x)$ is a Multilayer Perceptron (MLP). The structure of both the posterior network $q_\alpha(z|x, y_g)$ and the prior network $p_\beta(z|x)$ is a bidirectional Long Short-Term Memory (Bi-LSTM) [28] layer followed by an MLP. The Bi-LSTM is used to encode the tweet text. The only difference between the

Figure 3.3: The structure of the baseline CVAE model during training. Bi-LSTM is a bidirectional LSTM layer. MLP is a Multilayer Perceptron. $x$ is the embedded input text. $y_g$ is the ground truth hate group label. $\hat{y}_g$ is the output prediction of the hate group. $p$ is the posterior distribution of the latent variable $z$ while $p'$ is the prior distribution of $z$. Note that this structure is used for training. During testing, the posterior network is replaced by the prior network to compute $\hat{y}_g$ and thus $y_g$ is not available during testing.

posterior and the prior network is that for the posterior network the input for the MLP is the encoded text concatenated with the group label while for the prior network only the encoded text is fed forward. During testing, the posterior network will be replaced by the prior network to generate the distribution over the latent variable (i.e. $p'$ will replace $p$). Thus during testing, the ground truth labels will not be used to make predictions.

We assume the latent variable $z$ has a multivariate Gaussian distribution: $p = \mathcal{N}(\mu, \Sigma)$ for the posterior network, and $p' = \mathcal{N}(\mu', \Sigma')$ for the prior network. The detailed

computation process is as follows:

$$e = f(x) \tag{3.5}$$

$$\mu, \Sigma = s(y_g \oplus e) \tag{3.6}$$

$$\mu', \Sigma' = s'(f'(x)) \tag{3.7}$$

where $f$ is the Bi-LSTM function and $e$ is the output of the Bi-LSTM layer at the last time step. $s$ is the function of the MLP in the posterior network and $s'$ is that of the MLP in the prior network. The notation $\oplus$ means concatenation. The latent variables $z$ and $z'$ are randomly sampled from $\mathcal{N}(\mu, \Sigma)$ and $\mathcal{N}(\mu', \Sigma')$, respectively. During training, the input for the likelihood network is $z$:

$$\hat{y}_g = w(z) \tag{3.8}$$

where $w$ is the function of the MLP in the likelihood network. During testing, the prior network will substitute for the posterior network. Thus for testing, the input for the likelihood is $z'$ instead of $z$:

$$\hat{y}_g = w(z') \tag{3.9}$$

## 3.5 Hierarchical CVAE

One problem with the above CVAE method is that it utilizes the group label for training, but ignores the available hate group category (ideology) information of the hate speech. As mentioned in Section 3.1, hate groups can be divided into different categories in terms of ideologies. Each hate group belongs to a specific hate group category. Considering this hierarchical structure, the hate category information can

potentially help the model to better capture the subtle differences between the hate speech from different hate groups. Therefore, we extend the baseline CVAE model to incorporate the category information. In this case, the objective function is as follows:

$$
\begin{aligned}
Obj &= \sum_{(x,y_c,y_g)\in X} \log p(y_c, y_g|x) \\
&= \sum_{(x,y_c,y_g)\in X} \log p(y_c|x) + \log p(y_g|x, y_c)
\end{aligned}
\tag{3.10}
$$

where $y_c$ is the hate group category label and

$$
p(y_c|x) = \int_{z_c} p(y_c|z_c, x)p(z_c|x)dz_c
\tag{3.11}
$$

$$
p(y_g|x, y_c) = \int_{z_g} p(y_g|z_g, x, y_c)p(z_g|x, y_c)dz_g
\tag{3.12}
$$

where $z_c$ and $z_g$ are latent variables. Therefore, the ELBO of our method is the sum of the ELBOs of $\log p(y_c|x)$ and $\log p(y_g|x, y_c)$:

$$
\begin{aligned}
ELBO = &E[\log p(y_c|z_c, x)] - \\
&D_{KL}[q(z_c|x, y_c)||p(z_c|x)] + \\
&E[\log p(y_g|z_g, x, y_c)] - \\
&D_{KL}[q(z_g|x, y_c, y_g)||p(z_g|x, y_c)]
\end{aligned}
\tag{3.13}
$$

During testing, the prior networks will substitute the posterior networks and the ground truth labels $y_c$ and $y_g$ are not utilized. Hence the prior $p(z_g|x, y_c)$ in the above equation cannot be parametrized as a network that directly takes the ground truth label $y_c$ and $x$ as inputs. Instead, we parameterize it as shown in the right part of Figure 3.4. We assume that $u'$ is trained to be a latent representation of $y_c$, so we use $u'$ and $x$ as inputs for this prior network.

According to the ELBO above, the corresponding loss function $\mathcal{L}$ is the combination of the loss function for the category classification ($\mathcal{L}_c$) and that for the group classification ($\mathcal{L}_g$).

$$
\begin{aligned}
\mathcal{L} =& \mathcal{L}_c + \mathcal{L}_g, \text{ where} \\
\mathcal{L}_c =& E_{z_c \sim q_\alpha(z_c|x,y_c)}[-\log p_\varphi(y_c|z_c, x)] + \\
& D_{KL}[q_\alpha(z_c|x, y_c)||p_\beta(z_c|x)], \text{ and} \\
\mathcal{L}_g =& E_{z_g \sim q_\eta(z_g|x,y_c,y_g)}[-\log p_\theta(y_g|z_g, x, y_c)] + \\
& D_{KL}[q_\eta(z_g|x, y_c, y_g)||p_\gamma(z_g|x, u)]
\end{aligned}
\tag{3.14}
$$

Figure 3.4 shows the structure of our model for training. By assuming the latent variables $z_c$ and $z_g$ have multivariate Gaussian distributions, the actual outputs of the posterior and prior networks are the mean and variance: $p_c = \mathcal{N}(\mu_c, \Sigma_c)$, $p_g = \mathcal{N}(\mu_g, \Sigma_g)$ for the two posterior networks, and $p'_c = \mathcal{N}(\mu'_c, \Sigma'_c)$, $p'_g = \mathcal{N}(\mu'_g, \Sigma'_g)$ for the two prior networks. Note that in addition to these four distributions, there is another distribution $p_x = \mathcal{N}(\mu_x, \sigma_x)$ as shown in Figure 3.4. This distribution is generated only with the input text $x$, so it provides the basic distribution (in order to avoid confusion, we call it the basic distribution instead of the prior distribution) for both two posterior networks and two prior networks. With this basic distribution, the two posterior networks only need to capture the additional signals learned from $x$ and labels. Similarly, the prior networks only need to learn the additional signals learned by the posterior networks. The detailed computation process during training is shown as the following equations:

Figure 3.4: The structure of the HCVAE during training. $x$ is the embedded input text. $y_c$ and $y_g$ are the ground truth hate category label and hate group label. $\hat{y}_c$ and $\hat{y}_g$ are the output predictions of the hate category and hate group. $z_c$, $z_g$, and $z'_c$ are latent variables. In the left dotted box are two posterior networks. In the right dotted box are two prior networks. Note that this structure is used for training. During testing, the posterior networks will be substituted by the posterior networks (i.e. the left dotted box will be substituted with the right one) to compute $\hat{y}_c$ and $\hat{y}_g$. Thus $y_c$ and $y_g$ are not available during testing.

$$e = f(x) \tag{3.15}$$

$$\mu_x, \Sigma_x = s_x(e) \tag{3.16}$$

$$\mu_c, \Sigma_c = s_c(y_c \oplus e) \tag{3.17}$$

$$\mu_g, \Sigma_g = s_g(y_c \oplus y_g \oplus e) \tag{3.18}$$

$$\mu'_c, \Sigma'_c = s'_c(f'_c(x)) \tag{3.19}$$

where $f$ is the Bi-LSTM function and $e$ is the output of the Bi-LSTM layer at the last

time step. $s_x$, $s_c$, $s_g$, and $s'_c$ are the functions of four different MLPs. $f'_c$ is the Bi-LSTM function in the prior network $p_\beta(z_c|x)$. The latent variables $z_x$, $z_c$, $z_g$, and $z'_c$ are randomly sampled from the Gaussian distributions $\mathcal{N}(\mu_x, \Sigma_x)$, $\mathcal{N}(\mu_c, \Sigma_c)$, $\mathcal{N}(\mu_g, \Sigma_g)$, and $\mathcal{N}(\mu'_c, \Sigma'_c)$ respectively. As mentioned above, $p_x$ is the basic distribution while $p_c$, $p_g$, $p'_c$, and $p'_g$ are trained to capture the additional signals. Therefore, $z_x$ is added to $z_c$ and $z'_c$, then the results $u$ and $u'$ are further added to $z_g$ and $z'_g$, respectively, as shown in the following equations:

$$u' = z_x + z'_c \tag{3.20}$$

$$\mu'_g, \Sigma'_g = s'_g(u' \oplus f'_g(x)) \tag{3.21}$$

$$u = z_x + z_c \tag{3.22}$$

$$v = u + z_g \tag{3.23}$$

where $f'_g$ is the Bi-LSTM function and $s'_g$ is the function of the MLP in the prior network $p_\gamma(z_g|x, u)$. $+$ is element-wise addition. During training, $u$ and $v$ are fed into the likelihood networks:

$$\hat{y}_c = w_c(e \oplus u) \tag{3.24}$$

$$\hat{y}_g = w_g(e \oplus v) \tag{3.25}$$

where $w_c$ and $w_g$ are the functions of the MLPs in two likelihood networks. During testing, the prior networks will substitute the posterior networks, so the latent variable $z'_g$ is randomly sampled from the Gaussian distributions $\mathcal{N}(\mu'_g, \Sigma'_g)$, and the last four

equations above (Equation 3.22 - 3.25) will be replaced by the following ones:

$$\hat{y}_c = w_c(e \oplus u') \tag{3.26}$$

$$v' = u' + z_g' \tag{3.27}$$

$$\hat{y}_g = w_g(e \oplus v') \tag{3.28}$$

Algorithm 2 and 3 illustrate the complete training and testing process. $BCE$ refers to

---

**Algorithm 2** Training Algorithm

---

1: **function** TRAIN($X$)
2:     randomly initialize network parameters;
3:     **for** $epoch = 1, E$ **do**
4:         **for** $(text, category, group)$ in $X$ **do**
5:             get embeddings $x$ and one-hot vectors $y_c$, $y_g$;
6:             compute $e$ with the Bi-LSTM;
7:             compute $\mu_x$, $\Sigma_x$, $\mu_c$, $\Sigma_c$, $\mu_g$, $\Sigma_g$;
8:             sample $z_x = reparameterize(\mu_x, \Sigma_x)$;
9:             sample $z_c = reparameterize(\mu_c, \Sigma_c)$;
10:             sample $z_g = reparameterize(\mu_g, \Sigma_g)$;
11:             $u = z_x + z_c$, $v = u + z_g$;
12:             compute $\hat{y}_c$ and $\hat{y}_g$ according to Eq. 3.24- 3.25;
13:             $\mathcal{L}_{REC} = BCE(\hat{y}_c, y_c) + BCE(\hat{y}_g, y_g)$;
14:             compute $\mu_c'$, $\Sigma_c'$;
15:             sample $z_c' = reparameterize(\mu_c', \Sigma_c')$;
16:             $u' = z_x + z_c'$;
17:             compute $\mu_g'$, $\Sigma_g'$;
18:             $\mathcal{L}_{KL} = D_{KL}(p_c||p_c') + D_{KL}(p_g||p_g')$;
19:             $\mathcal{L} = \mathcal{L}_{KL} + \mathcal{L}_{REC}$;
20:             update network parameters on $\mathcal{L}$;
21:         **end for**
22:     **end for**
23: **end function**

---

the Binary Cross Entropy loss. Note that during training, the ground truth category labels and group labels are fed to the posterior network to generate latent variables. But during testing, the latent variables are generated by the prior network, which only utilizes the texts as inputs.

---

**Algorithm 3** Testing Algorithm

---

1: **function** TEST($X$)
2:   **for** *text* in $X$ **do**
3:     get embeddings $x$;
4:     compute $e$ with the Bi-LSTM;
5:     compute $\mu_x$, $\Sigma_x$, $\mu'_c$, $\Sigma'_c$;
6:     sample $z_x = reparameterize(\mu_x, \Sigma_x)$;
7:     sample $z'_c = reparameterize(\mu'_c, \Sigma'_c)$;
8:     $u' = z_x + z'_c$;
9:     compute $\mu'_g$, $\Sigma'_g$;
10:     sample $z'_g = reparameterize(\mu'_g, \Sigma'_g)$;
11:     $v' = u' + z'_g$;
12:     compute $\hat{y}_c$ and $\hat{y}_g$ according to Eq. 3.26- 3.28;
13:   **end for**
14: **end function**

---

# 3.6 Experiments and Discussion

## 3.6.1 Experimental Settings

In addition to the discriminative CVAE model described in Section 3.4, we implement other baseline methods and an upper bound model as follows.

**Support Vector Machine (SVM)**: We implement an SVM model with linear kernels. We use L2 regularization and the coefficient is 1. The input features are the Term Frequency Inverse Document Frequency (TF-IDF) values of up to 2-grams.

**Logistic Regression (LR)**: We implement the Logistic Regression model with L2 regularization. The penalty parameter C is set to 1. Similar to the SVM model, we use TF-IDF values of up to 2-grams as features.

**Character-level Convolutional Neural Network (Char-CNN)**: We implement the character-level CNN model for text classification as described in [44]. It is 9 layers deep with 6 convolutional layers and 3 fully-connected layers. The configurations of the convolutional layers are kept the same as those in [44].

**Bi-LSTM**: The model consists of a bidirectional LSTM layer followed by a linear layer. The embedded tweet text $x$ is fed into the LSTM layer and the output at the last time step is then fed into the linear layer to predict $\hat{y}_g$.

**Upper Bound**: The upper bound model also consists of a bidirectional LSTM layer followed by a linear layer. The difference between this model and Bi-LSTM is that it takes the tweet text $x$ along with the ground truth category label $y_c$ as input during both training and testing. The LSTM layer is used to encode $x$. The encoding result is concatenated with the ground truth category label and then fed into the linear layer to give the prediction of the hate group $\hat{y}_g$. Since it utilizes $y_c$ for prediction during testing, this model sets an upper bound performance for our method.

For the baseline CVAE, Bi-LSTM, the upper bound model, and the HCVAE, we use randomly initialized word embeddings with size 200. All the neural network models are optimized by the Adam optimizer with a learning rate 1e-4. The batch size is 20. The hidden size of all the LSTM layers in these models is 64 and all the MLPs are three-layered with non-linear activation function Tanh. For the CVAE and the HCVAE, the size of the latent variables is set to 256.

All the baseline models and our model are evaluated on two datasets. We first use the complete dataset for training and testing. 90% of the instances are used for training and the rest for testing. In order to evaluate the robustness of the baseline models and our model, we also randomly selected a subset (10%) of the original dataset for training while the testing dataset is fixed. Since the upper bound model is used to set an upper bound on performance, we do not evaluate it on the smaller training dataset.

| Dataset | Complete | | | Subset | | |
|---|---|---|---|---|---|---|
| Metric | Macro-F1 | Micro-F1 | Wei.-F1 | Macro-F1 | Micro-F1 | Wei.-F1 |
| upper bound | .697 | .904 | .881 | – | – | – |
| SVM + tf-idf | .653 | .834 | .835 | **.521** | .771 | .772 |
| LR + tf-idf | .586 | .787 | .792 | .494 | .727 | .736 |
| Char-CNN | .604 | .840 | .853 | .457 | .730 | .744 |
| Bi-LSTM | .627 | .767 | .745 | .353 | .599 | .570 |
| CVAE | .520 | .799 | .830 | .453 | .774 | .784 |
| HCVAE (our) | **.664** | **.844** | **.858** | .469 | **.787** | **.799** |

Table 3.1: Experimental results. Complete: The performance achieved when 90% of the entire dataset is used for training. Subset: The performance achieved when only 10% of the dataset is used for training. Wei.-F1: Weighted-F1. The best results are in bold.

## 3.6.2   Experimental Results

The experimental results are shown in Table 3.1. The testing dataset keeps the imbalanced distribution of the original data, so the Macro-F1 score of each method is significantly lower than the Micro-F1 score and the Weighted-F1 score. Comparing the performance of the Bi-LSTM model and that of the baseline CVAE model shows that although CVAE is traditionally used as a generative model, it can also achieve competitive performance as a discriminative model. All the methods achieve lower F1 scores when using the smaller training dataset. Due to the imbalanced distribution of our dataset, half of the 40 groups have less than 1k tweets in the smaller training dataset, which leads to a sharp decline in the Macro-F1 scores of all the methods. However, the performance of both CVAE-based models degrades less than that of the other two neural network models (the Bi-LSTM model and the Char-CNN model) according to the Micro-F1 Weighted-F1 scores. These two CVAE-based models tend to be more robust when the size of the training dataset is smaller. The difference between the CVAE-based models and the other two models is that both the Bi-LSTM model and the Char-CNN model directly compress the input text into a fixed-length latent variable while the CVAE

Figure 3.5: This figure compares a subset (17 hate groups of 6 categories) of the predictions made by the baseline CVAE and the HCVAE. The X and Y axes are the index of each hate group. The hate groups of the same category are grouped together as shown in the dashed squares. The color of the grid $(i, j)$ is mapped from $r_{hcvae}(i, j) - r_{cvae}(i, j)$, where $r(i, j)$ is the fraction of the group $i$'s instances among the instances predicted as the group $j$. A higher difference value corresponds to a lighter color. A grid darker than the background is mapped from a negative value while a grid lighter than the background is mapped from a positive one.

model explicitly models the posterior and likelihood distributions over the latent variable. The result is that, during testing, the inference strategy of the Bi-LSTM model and the Char-CNN model is actually comparing the compressed text to the compressed instances in the training dataset while the strategy of the CVAE-based models is to compare the prior distributions over the latent variable. Compared with the compressed text, prior distributions may capture higher-level semantic information and thus enable better generalization.

By further utilizing the hate category label for training, the HCVAE outperforms the baseline CVAE on all three metrics. Figure 3.5 illustrates the difference between the predictions made by the HCVAE and the CVAE. As shown in the figure, most of the lighter girds are in the dashed squares while most of the darker girds are outside.

F1 score



Figure 3.6: The F1 score achieved by our method on each hate group. The Y-axis is the F1 score. The X-axis is the hate group sorted by the number of instances in the dataset from larger to smaller. The dashed line shows the F1 scores and the solid line is the corresponding trendline.

This indicates that the HCVAE model tends to correct the prediction of the CVAE's misclassified instances to be closer to the ground truth. In most cases, the misclassified instances can be directly corrected to the ground truth (the diagonal). In other cases, they are not corrected to the ground truth but are corrected to the hate groups of the same categories as the ground truth (the dashed squares). This shows that additional ideology information is useful for the model to better capture the subtle differences among tweets posted by different hate groups.

Although our method outperforms the baseline methods, there is still a gap between its performance and the upper bound performance. We analyze this in the following section.

### 3.6.3    Error Analysis

As mentioned above, in some cases our model misclassified the tweet as a group under the same category as the ground truth. Take, for instance, the tweet from *Group1*: *The only good #muslim is a #dead #muslim.* Our HCVAE model predicts it as from *Group2*. But both *Group1* and *Group2* are under the category "neo nazi". One possible reason for

this kind of error is that the imbalance of the dataset adversely affects the performance of our method. Figure 3.6 shows the F1 scores achieved by the HCVAE on each group. The performance of the model tends to be much lower when the number of the group's training instances decreases. Although we use the weighted random sampler to alleviate the problem of the imbalanced dataset, repeating the very limited data instances (less than 3k) of a group during training cannot really help the posterior and prior networks to give a reasonable distribution that can generalize well on unseen inputs. Therefore, when the model comes into the instances in the testing dataset, the performance can be less satisfactory. This is a common problem for all the methods, which is the cause of the significantly lower Macro-F1 scores.

Another type of error is caused by the noisy dataset. Take, for instance, the tweet from *Group3* under the category "ku klux klan": *we must secure the existence of our people and future for the White Children !!* Our model classifies it as from *Group4* under the category "neo nazi", which makes sense but is an error.

## 3.7 Conclusion

In this chapter, we explore the CVAE as a discriminative model and find that it can achieve competitive results. In addition, the performance of the CVAE-based models tends to be more stable than that of the others when the dataset gets smaller. We further propose an extension of the basic discriminative CVAE model to incorporate the hate group category (ideology) information for training. Our proposed model has a hierarchical structure that mirrors the hierarchical structure of hate groups and ideologies. We apply the HCVAE to the task of fine-grained hate speech classification, but this Hierarchical CVAE framework can be easily applied to other tasks where the hierarchical structure of the data can be utilized.

41

# Chapter 4

# Lifelong Learning of Fine-grained Hate Speech Characterization

In the previous chapter, we tackle the fine-grained hate speech characterization task in the isolated learning setting, where the dataset is assumed to be static and the classifiers are trained on the given training dataset. However, the amount of data in social media increases every day, and the hot topics change rapidly, requiring the classifiers to be able to continuously adapt to new data and accumulate knowledge without forgetting what was previously learned. This ability, referred to as lifelong learning, is crucial for the real-world application of hate speech classifiers in social media. In this chapter, we propose lifelong learning of hate speech characterization on social media, and novel methods to alleviate catastrophic forgetting. Experimentally, we show that our method achieves better performance than the commonly-used lifelong learning techniques.

Figure 4.1: An illustration of our proposed task. $hg_i$: the $i$th hate group. The model is trained on a sequence of sub-datasets, split by their hate ideologies, e.g., anti-Muslim and Kuklux Klan. The task on each sub-dataset is to identify the hate group given the tweet.

## 4.1 Introduction

Fine-grained classification is necessary for fine-grained hate speech analysis. The Southern Poverty Law Center (SPLC) monitors hate groups throughout the United States by a variety of methodologies to determine the activities of groups and individuals, including reviewing hate group publications [59]. Therefore, in the previous chapter, we propose a more fine-grained hate speech classification task that attributes hate groups to individual tweets. However, a common limitation of all the research mentioned above is that they assume the dataset to be static and train the classifiers on each isolated dataset, i.e., isolated learning, ignoring the rapid increase of the amount of data in social media and the rapid change of the hot topic.

A report from L1ght [61], a company that specializes in measuring online toxicity, suggests that amid the growing threat of the coronavirus, there has been a 900% growth in hate speech towards China and Chinese people on Twitter since February 2020. As a result of the rapid change in social media content, hate speech classifiers are required to be able to continuously learn and accumulate knowledge from a stream of data, i.e., lifelong

learning. Learning on each portion of the data is considered as a task, so a stream of tasks is joined to be trained sequentially. In this work, we propose a novel lifelong fine-grained hate speech classification task, as illustrated in Figure 4.1. The models trained by isolated learning tend to face catastrophic forgetting [62, 63, 64, 65] due to a non-stationary data distribution in lifelong learning. To address this problem, an extensive body of work has been proposed for various lifelong learning tasks. However, our experiments show that the commonly-used lifelong learning methods still exhibit catastrophic forgetting in our proposed tasks. One important difference between the Twitter hate group dataset and the other image datasets commonly used in lifelong learning studies is that the similarity among the different tasks is unstable and relatively low, as indicated by the low average Jaccard Indexes of the topic words in Table 4.1. To alleviate this problem, we introduce Variational Representation Learning (VRL) to distill the knowledge from each task into a latent variable distribution. We also augment the model with a memory module and adapt the clustering algorithm, LB-SOINN (Load-Balancing Self-Organizing Incremental Neural Network), to select the most important samples from the training dataset of each task. Our implementation is publicly available[1].

Our contributions are three-fold:

- This is the first paper on lifelong learning of fine-grained hate speech classification.

- We propose a novel method that utilizes VRL along with an LB-SOINN memory module to alleviate catastrophic forgetting resulted from a severe change in data distribution.

- Experimental results show that our proposed method outperforms the state-of-the-art significantly on the average F1 scores.

---

[1]The implementation is available here: https://www.aclweb.org/anthology/2021.naacl-main.183/

| Ideology | Avg. JI | Keywords |
|---|---|---|
| Christian Identity | 0.019 | Jesus, Yahuwshua |
| Radical Tr. Catholic | 0.031 | catholic, remnant |
| Neo Confederate | 0.039 | southern, Free Dixie |
| Anti Semitism | 0.047 | Israel, Trump |
| Anti Catholic | 0.049 | Texe Marrs, truth |
| Hate Music | 0.049 | death, radio |
| Anti Muslim | 0.064 | Muslim, Islam |
| Black Separatist | 0.071 | black, panther |
| Racist Skinhead | 0.074 | shirt, white |
| Anti Immigration | 0.075 | immigration, border |
| Holocaust Identity | 0.078 | Jewish, Trump |
| Neo Nazi | 0.091 | Hitler, white |
| Kuklux Klan | 0.100 | ni**a, f**king |
| Anti LGBTQ | 0.100 | family, marriage |
| White Nationalist | 0.105 | white, America |

Table 4.1: Information about the 15 hate ideologies. Tr.: Traditional. Avg JI: the average of the Jaccard Index between the topic words of one ideology and those of another ideology. The topic words are extracted by Latent Dirichlet Allocation (LDA) [66]. The top 2 most frequent topic words are selected as keywords.

## 4.2    Related Work

Most research on lifelong learning alleviates catastrophic forgetting in the following three directions.

**Regularization-based Methods:** These methods impose constraints on the weight update. The goal of the constraints is to minimize deviation from trained weights when training on a new task. The constraints are generally modeled by additional regularization terms [67, 68, 69, 70, 71]. Elastic Weight Consolidation (EWC) [67] alleviates catastrophic forgetting by slowing down learning on the model parameters which are important to the previous task. The importance of the parameters is estimated by the Fisher information matrix. Instead of the Fisher information matrix, PathNet [69] uses agents embedded in the neural network to determine which parameters of the neural network can be reused for new tasks and the task-relevant pathways are frozen during

training on new tasks.

**Architecture-based Methods:** The main idea of this approach is to change architectural properties to dynamically accommodate new tasks, such as assigning a dedicated capacity inside a model for each task. Progressive Neural Networks [72] expand the model architecture by allocating a new column of neural network for each new task. In [73, 74], Convolutional Neural Network is combined with LB-SOINN for incremental online learning of object classes. Although they also use LB-SOINN in their work, the usage of LB-SOINN in this work is completely different. They use LB-SOINN to predict object class while our proposed method adapts the original LB-SOINN to calculate the importance of the training samples without making any prediction on the class. A problem with the methods in this category is that the available computational resources are limited in practice. As a result, the model expansion will be prohibited when the number of tasks increases to a certain degree.

**Data-based Methods:** These methods alleviate catastrophic forgetting by utilizing a memory module, which either stores a small number of real samples from previous tasks or distills knowledge from previous tasks. The main feature of Gradient Episodic Memory (GEM) [75] is the episodic memory, storing a subset of the samples from the observed tasks. GEM computes the losses on the episodic memories and treats them as inequality constraints, avoiding them to increase. Averaged GEM [76] is a more efficient version of GEM. Later work proposes a lifelong language learning model using a key-value memory module for sparse experience replay and local adaptation [77]. LAMOL [78] formulates lifelong language learning as a language modeling task and replays the generated pseudo-samples of previous tasks during training.

There are also studies combining multiple methods above. In [79], the architecture-based method is combined with the data-based method. In [80], the regularization method is combined with the data-based method for lifelong learning on relation ex-

traction. Our proposed method is also a combination of the regularization method and the data-based method but in a different way.

## 4.3   Data, Task, and Measures

We collect hate group data in the same way as in the previous chapter, where the tweet handles are collected based on the hate groups identified by SPLC. SPLC categorizes these hate groups according to their hate ideologies. For each hate ideology, the top three Twitter handles are selected in terms of the number of followers. The dataset includes all the content (tweets, retweets, and replies) posted with each Twitter account from the group's inception date, as early as 2009, until 2017. Altogether, the dataset consists of 42 hate groups from 15 different ideologies. Table 4.1 shows the 15 ideologies. Each instance in the dataset is a text tuple of (tweet, hate group name, hate ideology).

We separate the dataset by ideology. The reason is that various existing hate speech datasets collect data using keywords or hashtags [18, 10, 81], which have a strong relationship with hate ideologies or topics. We also observe that the hot spots of society can lead to a significant shift of major hate speech topics or the emergence of new hate ideologies on social media as mentioned in Section 4.1, indicating that the expansion of the hate speech dataset may be accompanied by the emergence of new hate ideologies.

Therefore, we separate the collected data into a sequence of 15 subsets according to their ideologies and sort them by the date of the first tweet post in each subset, from the earliest to the latest. The task on each subset is to identify the hate group given the tweet text. In the previous chapter, we propose a hierarchical Conditional Variational Autoencoder model for the fine-grained hate speech classification task. The architecture and the training process of their model require the number of classes to be pre-defined. However, we do not pre-define the number of classes in this task since such kind of

information is not available in the real-world application of lifelong learning. The model should be able to incorporate emerging hate groups at any time of training. In order to satisfy this condition, we formulate the task of identifying the group as a ranking task, instead of a classification task. For each tweet, we provide the model with a set of candidate groups, consisting of all the previously seen hate groups, including the ground truth group. The model takes each combination of the tweet and the candidate group as input and outputs a score. The corresponding loss function is:

$$\mathcal{L}_r = \sum_{(x,y_s)\in D} \sum_{y_i \in Y\setminus\{y_s\}} h(f_\theta(x, y_s) - f_\theta(x, y_i)) \tag{4.1}$$

where $x$ is the tweet text, and $y_s$ is the ground truth group of $x$. $Y$ is the candidate group set of $x$, which consists of all the seen hate groups until $x$ is observed by the model, including the ground truth group $y_s$ of $x$, so $y_i \in Y\setminus\{y_s\}$ is the negative candidate group of $x$. $f_\theta$ is the scoring model parameterized by $\theta$. $h(a) = \max(0, m - a)$, $m$ is the chosen margin.

Same as in other lifelong learning studies, we consider learning on each of the hate ideologies in the sequence as a task, so we have a sequence of 15 tasks. As mentioned in Section 4.1, the similarity among our tasks is unstable and relatively low. Therefore, when the model is continuously trained on the tasks, it may encounter a sudden change in vocabulary, topic, and input data distribution. This makes our tasks more challenging compared to the other lifelong learning tasks because the abrupt change can make the catastrophic forgetting problem more severe. This is also the reason that some techniques achieving significant improvement in the image classification tasks do not perform well on our task (see Section 4.5).

We use the average macro F1 score and average micro F1 score for evaluation.

$$\text{Average F1: } AvgF1(t) = \frac{1}{t} \sum_{i=1}^{t} F1_{t,i} \tag{4.2}$$

where $F1_{t,i}$ is the F1 score, either macro F1 or micro F1, achieved by the model on the $i$th task after being trained on the $t$th task. The larger this metric, the better the model.

## 4.4   Methods

As mentioned in Section 4.2, one way to alleviate catastrophic forgetting is to use a memory module, storing a small number of real samples from previous tasks, and a simple way to utilize the memorized samples is to replay the memory when training on a new task, such as mixing them with the training samples from the current task. The idea behind this approach is that the memorized samples should reflect the data distribution so that the replay of the memory can help the model make invariant predictions on the samples of the previous tasks. However, this approach may not work well when the size of the memory is small. The reason is that when there is only a small amount of data memorized, the memory is not able to reflect the data distribution of the previous task and thus the model can easily overfit on the memorized samples instead of generalizing to all the samples in the previous task.

We address this problem from two aspects. First, since the memory size is limited, it is beneficial to select the most representative training samples in the previous tasks to memorize. Second, simply storing the real training samples in the memory may not be sufficient to represent the knowledge of the previous tasks, so we need a better way to distill knowledge from the observed samples along with a method to utilize it when training on a new task. We combine two techniques: Variational Representation Learning

(VRL) and Load-Balancing Self-Organizing Incremental Neural Network (LB-SOINN) to achieve these goals. We propose a supervised version of LB-SOINN to select the most important training samples in the current task. VRL not only distills the knowledge from the current training task but also provides an appropriate hidden representation as input for the LB-SOINN, so we introduce VRL first.

### 4.4.1  Variational Representation Learning

The distilled knowledge of previous tasks can take various forms, but the key point is that it should be related to the data distribution of the corresponding task so that it can be utilized to alleviate catastrophic forgetting. Inspired by the Variational Autoencoder (VAE) [54], we consider the distribution of the hidden representation of the input data as the distilled knowledge.

The original VAE model is proposed for data generation, so the objective of the original VAE is:

$$Obj = \sum_{x \in X} \log p(x) \tag{4.3}$$

$$p(x) = \int_z p(x|z)p(z)dz \tag{4.4}$$

$z$ is the latent variable, i.e., the hidden representation of the input. Since the integration over $z$ is intractable, we instead try to maximize the corresponding evidence lower bound (ELBO) and the corresponding loss function is as follows:

$$\mathcal{L}_{vae} = \sum_{x \in X} E_{z \sim p_\alpha(z|x)}[-\log p_\varphi(x|z)]+ \\ D_{KL}[q_\alpha(z|x)||p_\beta(z)] \tag{4.5}$$

$p(x|z)$, $q(z|x)$, and $p(z)$ are the likelihood distribution, posterior distribution, and prior

distribution. $\alpha,\varphi$, and $\beta$ indicate parameterization. The loss function can be separated into two parts. The first part $E[-\log p(x|z)]$ is the reconstruction loss, trying to reconstruct the input text from the latent variable. It pushes $z$ to reserve as much information of the input as possible. This is consistent with our goal to learn knowledge of data distribution. The second part is $D_{KL}[q(z|x)||p(z)]$, where $D_{KL}$ is the Kullback–Leibler (KL) divergence. Minimizing it pushes the posterior and the prior distributions to be close to each other. By assuming the posterior $p(z|x)$ to be a multivariate Gaussian distribution $\mathcal{N}(\mu_z, \Sigma_z)$, the latent variable $z$ is sampled from $\mathcal{N}(\mu_z, \Sigma_z)$.

In the original VAE, $p(z)$ is chosen to be a simple Gaussian distribution $\mathcal{N}(0, 1)$. However, this is over-simplified in our task because different from the unsupervised generation task of the original VAE, our ranking task is supervised. Our task not only requires $z$ to contain information of the tweet text itself but also requires it to indicate the group information of the tweet. In other words, the distilled distribution should be conditioned on both the tweet and its group label to reflect the data distribution in a supervised task. Setting the prior to be the same for all the hate groups pushes $z$ or the distribution of $z$ to ignore the label information. Instead, the prior should be different for each hate group, so we replace $p(z)$ with $p(u|y_s)$, where $y_s$ is the group label of $x$ and $u$ is the latent variable. $p(u|y_s)$ is assumed to be a multivariate Gaussian distribution $\mathcal{N}(\mu_u, \Sigma_u)$. Note that the replacement itself can not guarantee $p(u|y_s)$ to be different for each hate group because the loss function in Equation 4.5 does not push $p(u|y_s)$ to satisfy this condition. However, the ranking loss function 4.1 fills in the gap. Therefore,

Figure 4.2: An illustration of our method. The dotted arrows indicate the computation of the loss. The light-colored dashed arrows illustrate the update of the memory module. Note that the layers in the rounded rectangle share parameter weight. There is only one encoder for the group input, followed by two linear layers. We make a copy of it in the figure for a clear illustration of loss computation. $\hat{x}$: the reconstructed tweet input. $s_1$, $s_2$: scores of $(x, y_s)$ and $(x, y_i)$ separately. $\mu_z^*$ and $\Sigma_z^*$ are the previously memorized distribution on the latent variable of $x$. $L_{rec}$ is the reconstruction loss, which is the first term in Equation 4.5.

our loss function on the current training task is a combination of these two.

$$\mathcal{L}_{cur} = \sum_{(x,y_s)\in D} \sum_{y_i\in Y\setminus\{y_s\}} h(f_\theta(x, y_s) - f_\theta(x, y_i))$$
$$+ E_{z\sim p_\alpha(z|x)}[-\log p_\varphi(x|z)]$$
$$+ D_{KL}[q_\alpha(z|x)||p_\beta(u|y_s)]$$

$$(4.6)$$

The right part of Figure 4.2 illustrates the computation process of VRL.

## 4.4.2   LB-SOINN Memory Module

VRL provides a way to summarize knowledge into latent variable distributions. However, we still need a method to utilize the learned distribution to alleviate catastrophic

forgetting. We do this by incorporating a memory module $D_{mem}$ to store a small subset of important training samples along with their latent variable distributions, so each sample stored in the memory is a tuple of $(x, y_z, q_{\alpha'}(z|x))$. Here $q_{\alpha'}(z|x)$ is the distribution computed when the model completes training on the task that $(x, y_z)$ belongs to. The memorized samples are taken as anchor points when training on a new task. We introduce a memory KL divergence loss to push $q_\alpha(z|x)$ computed when training on a new task to be close to the memorized distribution $q_{\alpha'}(z|x))$. Therefore, the complete loss function is:

$$
\begin{aligned}
\mathcal{L} &= \mathcal{L}_{cur} + D_{KLmem} \\
&= \mathcal{L}_{cur} + \sum_{(x,y_s)\in D_{mem}} D_{KL}[q_\alpha(z|x)||q_{\alpha'}(z|x))]
\end{aligned}
\tag{4.7}
$$

Since the size of the memory is limited, we introduce a supervised version of LB-SOINN to select the most important training samples in the current task. The input for the LB-SOINN is the hidden representation of the tweet text, which is $z$ in the case of Variational Representation Learning (see Figure 4.2). We refer readers to the original work [82] for a detailed explanation of LB-SOINN. The original LB-SOINN is an unsupervised clustering algorithm that clusters unlabeled data by topology learning. We utilize the topology learning of LB-SOINN instead of clustering since our task is supervised. Therefore, we make the following adjustments to the original LB-SOINN.

1) The criteria to add a new node: Add a new node to the node set if one of the following conditions is satisfied: a) The distance between the input and the winner is larger than the winner's threshold. b) The distance between the input and the second winner is larger than the second winner's threshold. c) The label of the input sample is not the same as the label of the winner.

2) Build connections between nodes: Connect the two nodes with an edge only if the winner and the second winner belong to the same class.

3) We disable the removal of edges whose ages are greater than a predefined parameter. We disable the deleting of nodes and the algorithm of updating the subclass labels of every node. The node label is the label of the instances assigned to it. Our adjusted algorithm guarantees that each node will only be assigned the samples from one class.

LB-SOINN keeps track of the density of each node, which is defined as the mean accumulated points of a node. A node gets points when there is an input sample assigned to it. If the mean distance of the node from its neighbors is large, we give low points to the node. In contrast, if the mean distance of the node from its neighbors is small, we give high points to the node. Therefore, the density of the node reflects the number of nodes close to it and also the number of samples assigned to it. We take the density of the node as a measurement of the importance of the samples assigned to the node. After the LB-SOINN finishes training on the samples from the current task, we sort the samples according to the density of the node they are assigned to and the top $K$ samples are selected to write to the memory. We divide the memory equally for each of the previous tasks, so $K = M/t$, where $M$ is the total memory size and $t$ is the number of observed tasks, including the current task. The old memory consists of samples from the previous $t - 1$ tasks and each task keeps $M/(t - 1)$ samples in the old memory. For each of the $t - 1$ tasks, the $M/(t - 1) - M/t$ samples with the lowest node densities are deleted, resulting in $K$ empty slots in the memory, which are then rewritten by the selected $K$ samples in the current task.

## 4.5 Experiments and Discussion

### 4.5.1 Experimental Settings

For each task, we randomly sample 5000 tweets from 80% of the collected data for training, 10% of the collected data for testing, and the rest 10% for development. We allow the model to make more than one pass over the training samples in the current task or the current memory during training. We compare our methods with the following methods:

**Fine-tuning:** The model contains two bidirectional LSTM encoders [28, 35, 36] to encode the tweet and the group separately. The score of the group is calculated as the cosine distance between the hidden state of the tweet encoder and that of the group encoder. This model is also the backbone model of all the methods described below, except Fine-tuning + BERT. The model is directly fine-tuned on the stream of tasks, one after another, by the ranking loss function in 4.1.

**Fine-tuning+BERT:** The training framework is the same as above, but each encoder is replaced by a pre-trained BERT model [15] followed by a linear layer. The linear layers are fine-tuned during training.

**Fine-tuning+RMR (Random Memory Replay):** We augment the fine-tuning method with an additional memory module. Same as in Section 4.4.2, the memory is divided equally for each task, but instead of using LB-SOINN, the $K$ samples are randomly sampled from the current training data and then rewrite $K$ random slots in the old memory.

**EWC:** EWC is a regularization-based method, adding a penalty term $\sum_i \frac{\lambda}{2} F_i (\theta_i - \theta_i^*)^2$ to the ranking loss function 4.1. $F_i$ is the diagonal of the Fisher information matrix $F$, $\theta$ is the model parameter, and $i$ labels each parameter. $\theta^*$ is the model parameter when the model finishes training on the previous task. $\lambda$ is set to 2e6 in our experiments.

**GEM:** We use the episodic memory in the original paper: the memory is populated with $m$ random samples from each task. $m$ is a predefined size of the episodic memory. We set $m = 100$ in our experiments, so each task can add 100 tweets to the memory. By the end of the 15 tasks, the total memory of GEM contains 1500 tweets.

**Multitask Learning:** The tasks are trained simultaneously. We mix the training data from multiple tasks to train the model. This setting does not follow the lifelong learning setting where the tasks are trained sequentially. We add this setting in our experiments to show the potential room for improvement concerning each lifelong learning method.

We do not compare our method with Support Vector Machine [40] or Logistic Regression, because they require the number of classes to be fixed and to be known in advance, which is unrealistic in our tasks. We also do not experiment with the HCVAE method introduced in the previous chapter since it also has this requirement, as mentioned in Section 4.3.

In all our experiments, we use 1-layer bi-LSTM as encoders except for the fine-tuning + BERT setting and we use cosine distance to measure similarity. The input of the group encoder is the concatenation of the group name and its hate ideology. We use 1-layer bidirectional GRU [83] as the decoder in VRL. The hidden size of the encoders and the decoders is 64. The latent variable size in VRL is 128. We use 300-dimensional randomly initialized word embeddings. All the neural networks are optimized by Adam optimizer with the learning rate 1e-4. The batch size is 64. The loss margin $m = 0.5$. The maximum number of training epochs for each task is set to 20. For LB-SOINN, $\lambda$=1000, $\eta$=1.04. The memory size is limited to 1000 tweets for all the methods using a memory module except GEM. We do not set episodic memory size for each task as GEM because for lifelong hate speech classification, the number of tasks keeps increasing in the real world, and assuming unlimited total memory is unrealistic.

| Number of observed tasks | t=5 | | t=10 | | t=15 | |
|---|---|---|---|---|---|---|
| Avg F1 score (%) | Macro | Micro | Macro | Micro | Macro | Micro |
| Multitask | 15.26 | 67.07 | 5.05 | 37.20 | 3.57 | 38.61 |
| Fine-tuning | 6.02 | 16.44 | 4.35 | 5.77 | 3.96 | 6.18 |
| Fine-tuning + BERT | 6.02 | 16.44 | 4.06 | 5.45 | 3.03 | 5.80 |
| Fine-tuning + RMR | 11.15 | 44.40 | 2.56 | 15.77 | 3.51 | 15.19 |
| EWC | 8.57 | 20.42 | 2.42 | 6.81 | 1.95 | 7.27 |
| GEM | **13.04** | 30.95 | 3.07 | 12.51 | 2.70 | 15.07 |
| Ours | 12.61 | **49.75** | **6.96** | **47.30** | **5.13** | **44.62** |

Table 4.2: Experimental results. RMR: random memory replay. The best results are in bold.

## 4.5.2    Experimental Results

The experimental results are shown in Table 4.2. We report the performance of each method after the model finishes training on the first 5 tasks, the first 10 tasks, and all the 15 tasks. The average macro-F1 score is much lower than the average micro-F1 score due to the imbalanced data of each task. The large performance gap between the multitask training and fine-tuning shows that there exists severe catastrophic forgetting and that the low average F1 scores in the fine-tuning setting are not due to the model capacity. Replacing the bi-LSTM encoder with the pre-trained BERT encoder does not improve the performance. This reconfirms that the low scores result from catastrophic forgetting, not model capacity. Actually, fine-tuning and fine-tuning with BERT achieves the same average F1 scores at $t = 5$ because both models completely forget the previous tasks after converging on the fifth task, so both models achieve the same F1 scores on the testing data of the fifth task while achieving 0 scores on the previous four tasks. Due to the large model capacity of BERT, fine-tuning with BERT tends to overfit on the training data more seriously, leading to a slight performance decline at $t = 10$ and $t = 15$ compared to using bi-LSTM encoders. Since model capacity is not the key factor to solve catastrophic forgetting, we simply use bi-LSTM as encoders in our model instead of BERT, considering the computational cost.

| Number of observed tasks | t=5 | | t=10 | | t=15 | |
|---|---|---|---|---|---|---|
| Avg F1 score (%) | Macro | Micro | Macro | Micro | Macro | Micro |
| Full Model | 12.61 | 49.75 | **6.96** | **47.30** | 5.13 | **44.62** |
| w/o $D_{KLmem}$ | **15.00** | **58.64** | 4.21 | 36.36 | 3.72 | 40.87 |
| w/o VRL | 11.05 | 35.03 | 4.53 | 13.69 | 3.65 | 11.28 |
| w/o LB-SOINN | 13.01 | 50.99 | 6.15 | 44.42 | **5.59** | 30.91 |

Table 4.3: Ablation study. w/o $D_{KLmem}$: $D_{KLmem}$ in Equation 4.7 is removed. w/o VRL: VRL is replaced by the model used in the fine-tuning setting, i.e., fine-tuning + LB-SOINN memory replay. w/o LB-SOINN: LB-SOINN memory replay is replaced by random memory replay, i.e., VRL + RMR. The best results are in bold.

Adding RMR to the fine-tuning setting achieves significant performance improvement, even better than EWC or GEM. This is related to the characteristic of our tasks mentioned in Section 4.3. EWC remembers previous tasks by slowing down the update of the model parameters important to them, which is more suitable for the sequence of tasks that are similar to each other. However, significant changes in vocabulary, topic, or input data distribution are very common in our sequence of tasks, making memory replay more efficient than EWC. The performance of GEM during the second half of the training is close to that of fine-tuning with RMR, but there exists a gap in the first half. The reason is that GEM sets an episodic memory for each task, of which the size is 100 in our experiments, so before the 10th task in the sequence, the size of the total memory available for GEM is less than that of the memory module used in the fine-tuning with RMR setting.

Although RMR improves performance, the average F1 scores still drop quickly when the number of tasks increases. In the late stage of sequential training, each task can only keep dozens of samples in the memory and the model is not able to generalize well based on the memory. Our method solves this problem by combining VRL and LB-SOINN memory replay. The performance of our model is better and more stable than the other methods when the number of tasks increases. Our method achieves higher scores than

multitask training in the last four columns of Table 4.3 because learning on one task is easier than learning on a mix of tasks simultaneously. Every model in our sequential training experiments can easily achieve high F1 scores on the current task, making a large contribution to the average F1 scores. However, when doing multitask training, the model loses this benefit.

To investigate the effect of our method, we conduct the ablation study as shown in Table 4.3. Removing $D_{KLmem}$ from the final loss function in Equation 4.7 does not lower the performance when the number of observed tasks is small ($t$=5) because each task can store a few hundreds of samples in the memory at the early stage of sequential training, which is sufficient for the model to learn the previous tasks. However, when the number of tasks increases, $D_{KLmem}$ shows its effect on alleviating catastrophic forgetting.

Fine-tuning+LB-SOINN (Table 4.3) does not perform as well as fine-tuning+RMR (Table 4.2), while VRL+LB-SOINN (i.e., full model) performs better than VRL+RMR (Table 4.3). The reason lies in the input for LB-SOINN. Compared to the hidden representations spread evenly in the hidden space, the hidden representations which are well-organized in different group clusters make it easier for LB-SOINN to learn a reasonable topology structure of the training samples. VRL achieves this by explicitly pushing the hidden representation of tweets to follow a learned multivariate Gaussian distribution unique to each group. On the other hand, directly using the hidden state of the tweet encoder does not exhibit such kind of characteristics. VRL not only distills task knowledge but also provides an appropriate input for LB-SOINN, as stated in Section 4.4.

### 4.5.3   Error Analysis

Although our model achieves significant improvement over the baseline methods, we

Figure 4.3: The testing results of the first 5 tasks in the sequence when our model is trained on the first 10 tasks.

observe that our method does not perform well on the first task. As shown in Figure 4.3, there exists a large gap between the performance on the first task and the other tasks, and the micro-F1 score on the first task quickly drops to almost 0 when the number of observed tasks increases. We find the same results after we change the order of tasks in the sequence, so this is not the result of the task difficulty but is the result of our method. We find this problem is due to the reconstruction loss, which is the first part in Equation 4.5. The model observes a very limited number of tweets when training on the first task, making it difficult to learn the language model and reconstruct the tweet. As a result, the tweet representation learned on the first task may not contain the information we require, resulting in a large performance gap. When the number of observed tasks increases, this problem goes away quickly. We anticipate that pre-training the VAE in our model (the left branch in Figure 4.2) on a large Twitter corpus can alleviate this problem at the beginning of training.

## 4.6 Conclusion

In this paper, we introduce the lifelong hate speech classification task and propose

to use the VRL and LB-SOINN memory module to alleviate catastrophic forgetting. Our proposed method has the potential to benefit other lifelong learning tasks where the similarity between the contiguous tasks can be low. We make our implementation freely available to facilitate more application and investigation of our method in the future.

# Chapter 5

# Hate Symbol Deciphering

In the previous chapters, the computational models we discussed to understand hate speech frame the problem as a classification task, ignoring the understanding of hate symbols (e.g., *14 words*, *kigy*) and their secret connotations. In this chapter, we propose a novel task of deciphering hate symbols. To do this, we leverage the Urban Dictionary and collected a new, symbol-rich Twitter corpus of hate speech. We investigate neural network latent context models for deciphering hate symbols. More specifically, we study Sequence-to-Sequence models and show how they are able to crack the ciphers based on context. Furthermore, we propose a novel Variational Decipher and show how it can generalize better to unseen hate symbols in a more challenging testing setting.

## 5.1 Introduction

The vast quantity of hate speech on social media can be analyzed to study online abuse. In recent years, there has been a growing trend of developing computational models of hate speech, the majority of which focus solely on modeling hate speech as a binary or multiclass classification task [1, 29, 18, 10].

Figure 5.1: An example tweet with hate symbols.

While developing new features for hate speech detection certainly has merits, we believe that understanding hate speech requires us to design computational models that can decipher hate symbols that are commonly used by hate groups. Figure 5.1 shows an example usage of hate symbols in an otherwise seemingly harmless tweet that promotes hate. For example, *Aryan Warrior* is a longstanding racist prison gang based in the Nevada prison system. *WPWW* is the acronym for *White Pride World Wide*. The hate symbols *1488* and *2316* are more implicit. *14* symbolizes the 14 words: *"WE MUST SECURE THE EXISTENCE OF OUR PEOPLE AND A FUTURE FOR WHITE CHILDREN"*, spoken by members of the Order neo-Nazi movement. *H* is the 8th letter of the alphabet, so *88=HH=Heil Hitler*. Similarly, *W* is the 23rd and *P* is the 16th letter of the alphabet, so *2316=WP=White Power*.

In this work, we propose the first models for deciphering hate symbols. We investigate two families of neural network approaches: the Sequence-to-Sequence models [84, 83] and a novel Variational Decipher based on the Conditional Variational Autoencoders [54, 43, 42]. We show how these neural network models are able to guess the meaning of hate symbols based on context embeddings and even generalize to unseen hate symbols during testing. Our contributions are three-fold:

- We propose a novel task of learning to decipher hate symbols, which moves beyond the standard formulation of hate speech classification settings.

- We introduce a new, symbol-rich tweet dataset for developing computational models of hate speech analysis, leveraging the Urban Dictionary.

- We investigate a sequence-to-sequence neural network model and show how it is able to encode context and crack hate symbols. We also introduce a novel Variational Decipher, which generalizes better in a more challenging setting.

## 5.2  Related Work

### 5.2.1  Text Normalization in Social Media

The proposed task is related to text normalization focusing on the problems presented by user-generated content in online sources, such as misspelling, rapidly changing out-of-vocabulary slang, short-forms and acronyms, punctuation errors or omissions, etc. These problems usually appear as out-of-vocabulary words. Extensive research has focused on this task [85, 86, 87, 88]. However, our task is different from the general text normalization in social media in that instead of the out-of-vocabulary words, we focus on the symbols conveying hateful meaning. These hate symbols can go beyond lexical variants of the vocabulary and thus are more challenging to understand.

### 5.2.2  Machine Translation

An extensive body of work has been dedicated to machine translation. In [89], the authors study a number of natural language decipherment problems using unsupervised learning. Later work further frames the task of machine translation as decipherment

and tackles it without parallel training data [90]. Machine translation using deep learning (Neural Machine Translation) has been proposed in recent years and Sequence-to-Sequence (Seq2Seq) learning with Recurrent Neural Networks (RNN) has been successfully applied to this field [83, 84]. The translation performance is further improved using the attention mechanism [91]. Google's Neural Machine Translation System (GNMT) employs a deep attentional LSTM network with residual connections [92]. Recently, machine translation techniques have been also applied to explain non-standard English expressions [93]. However, our deciphering task is not the same as machine translation in that hate symbols are short and cannot be modeled as language.

Our task is more closely related to [94] and [95]. In [94], the authors propose using neural language embedding models to map the dictionary definitions to the word representations, which is the inverse of our task. In [95], the authors propose the definition modeling task. However, in their task, for each word to be defined, its pre-trained word embedding is required as an input, which is actually the prior knowledge of the words. However, such kind of prior knowledge is not available in our deciphering task. Therefore, our task is more challenging and is not simply a definition modeling task.

### 5.2.3   Hate Speech Analysis

Closely related to our work are [96, 97]. In [96], the proposed RNN model supplemented by an attention mechanism can highlight suspicious words for free, without including highlighted words in the training data. In [97], the authors propose a weakly-supervised approach that jointly trains a slur learner and a hate speech classifier. While their work contributes to the automation of harmful content detection and the highlighting of suspicious words, our work builds upon these contributions by providing a learning mechanism that deciphers suspicious hate symbols used by communities of hate to bypass

automated content moderation systems.

## 5.3   Data and Task

We first collect hate symbols and the corresponding definitions from the Urban Dictionary. Each term with one of the following hashtags: *#hate, #racism, #racist, #sexism, #sexist, #nazi* is selected as a candidate and added to the set $S_0$. We collected a total of 1,590 terms. Next, we expand this set by different surface forms using the Urban Dictionary API. For each term $s_i$ in set $S_0$, we obtain a set of terms $R_i$ that have the same meaning as $s_i$ but with different surface forms. For example, for the term *brown shirt*, there are four terms with different surface forms: *brown shirt, brown shirts, Brownshirts, brownshirt*. Each term in $R_i$ has its own definition in Urban Dictionary, but since these terms have exactly the same meaning, we select a definition $d_i$ with the maximum upvote/downvote ratio for all the terms in $R_i$. For example, for each term in the set $R_i=\{$*brown shirt, brown shirts, Brownshirts, brownshirt*$\}$, the corresponding definition is *"Soldiers in Hitler's storm trooper army, SA during the Nazi regime..."* After expanding, we obtain 2,105 distinct hate symbol terms and their corresponding definitions. On average, each symbol consists of 9.9 characters, 1.5 words. Each definition consists of 96.8 characters, 17.0 words.

For each of the hate symbols, we collect all tweets from 2011-01-01 to 2017-12-31 that contain exactly the same surface form of hate symbol in the text. Since we only focus on hate speech, we train an SVM [40] classifier to filter the collected tweets. The SVM model is trained on the dataset published in [18]. Their original dataset contains three labels: Sexism, Racism, and None. Since the SVM model is used to filter non-hate speech, we merge the instances labeled as sexism and racism, then train the SVM model to do binary classification. After filtering out all the tweets classified as non-hate, our

final dataset consists of 18,667 (tweet, hate symbol, definition) tuples.

We formulate hate symbol deciphering as the following equation:

$$Obj = \sum_{(u,s,d^*)\in X} \log p(d^*|(u,s)) \tag{5.1}$$

$X$ is the dataset, $(u, s, d^*)$ is the (tweet, symbol, definition) tuple in the dataset. The inputs are the tweet and the hate symbol in this tweet. The output is the definition of the symbol. Our objective is to maximize the probability of the definition given the (tweet, symbol) pair.

## 5.4   Methods

This objective function (Eq. 5.1) is very similar to that of machine translation. So we first try to tackle it based on the Sequence-to-Sequence model, which is commonly used in machine translation.

### 5.4.1   Sequence-to-Sequence Model

We implement an RNN Encoder-Decoder model with an attention mechanism based on [91]. We use GRU [83] for decoding. However, instead of also using GRU for encoding, we found that LSTM [28] performs better on our task. Therefore, our Seq2Seq model uses LSTM encoders and GRU decoders. An overview of our Seq2Seq model is shown in

Figure 5.2: Our Seq2Seq model. $u$, and $s_w$ are the word embeddings of the tweet text and hate symbol. $s_c$ is the character embedding of the symbol. $c_u$ is the encoded tweet and $h$ is the concatenated hidden states. $d$ is the generated text.

Figure 5.2. The computation process is shown as the following equations:

$$c_u, h_u = f_u(u) \tag{5.2}$$

$$c_{sw}, h_{sw} = f_{sw}(s_w) \tag{5.3}$$

$$c_{sc}, h_{sc} = f_{sc}(s_c) \tag{5.4}$$

$u$ is the word embedding of the tweet text, $s_w$ is the word embedding of the hate symbol, $s_c$ is the character embedding of the symbol. $f_u$, $f_{sw}$, and $f_{sc}$ are LSTM functions. $c_u$, $c_{sw}$, $c_{sc}$ are the outputs of the LSTMs at the last time step and $h_u$, $h_{sw}$, $h_{sc}$ are the hidden states of the LSTMs at all time steps. We use two RNN encoders to encode the symbol, one encodes at the word level and the other one encodes at the character level. The character-level encoded hate symbol is used to provide the feature of the surface form of the hate symbol while the word-level encoded hate symbol is used to provide the semantic information of the hate symbol. The hidden states of the two RNN encoders

for hate symbols are concatenated:

$$h = h_{sw} \oplus h_{sc} \tag{5.5}$$

$c_u$ is the vector of encoded tweet text. The tweet text is the context of the hate symbol, which provides additional information during decoding. Therefore, the encoded tweet text is also fed into the RNN decoder. The detailed attention mechanism and decoding process at time step $t$ are as follows:

$$w_t = \sigma(l_w(d_{t-1} \oplus e_{t-1})) \tag{5.6}$$

$$a_t = \sum_{i=1}^{T} w_{ti} h_i \tag{5.7}$$

$$b_t = \sigma(l_c(d_{t-1} \oplus a_t)) \tag{5.8}$$

$$o_t, e_t = k(c_u \oplus b_t, e_{t-1}) \tag{5.9}$$

$$p(d_t|u, s) = \sigma(l_o(o_t)) \tag{5.10}$$

$w_t$ is the attention weights at time step $t$ and $w_{ti}$ is the $i_{th}$ weight of $w_t$. $d_{t-1}$ is the generated word at last time step and $e_{t-1}$ is the hidden state of the decoder at last time step. $h_i$ is the $i_{th}$ time step segment of $h$. $l_w$, $l_c$, and $l_o$ are linear functions. $\sigma$ is a nonlinear activation function. $k$ is the GRU function. $o_t$ is the output and $e_t$ is the hidden state of the GRU. $p(d_t|u, s)$ is the probability distribution of the vocabulary at time step $t$. The attention weights $w_t$ are computed based on the decoder's hidden state and the generated word at time step $t-1$. Then the computed weights are applied to the concatenated hidden states $h$ of encoders. The result $a_t$ is the context vector for the decoder at time step $t$. The context vector and the last generated word are combined by a linear function $l_c$ followed by a nonlinear activation function. The result $b_t$ is concatenated with the

encoded tweet context $c_u$, and then fed into GRU together with the decoder's last hidden state $e_{t-1}$. Finally, the probability of each vocabulary word is computed from $o_t$.

### 5.4.2   Variational Decipher

The Variational Decipher is based on the CVAE model, which is another model that can be used to parametrize the conditional probability $p(d^*|(u,s))$ in the objective function (Eq. 5.1). Unlike the Seq2Seq model, which directly parametrizes $p(d^*|(u,s))$, our variational decipher formulates the task as follows:

$$
\begin{aligned}
Obj &= \sum_{(u,s,d^*)\in X} \log p(d^*|(u,s)) \\
&= \sum_{(u,s,d^*)\in X} \log \int_z p(d^*|z)p(z|(u,s))dz
\end{aligned}
\tag{5.11}
$$

where $z$ is the latent variable. $p(d^*|(u,s))$ is written as the marginalization of the product of two terms over the latent space. Since the integration over $z$ is intractable, we instead try to maximize the evidence lower bound (ELBO). Our variational lower bound objective is in the following form:

$$
\begin{aligned}
Obj =& E[\log p_\varphi(d^*|z,u,s)]- \\
& D_{KL}[p_\alpha(z|d^*,u,s)||p_\beta(z|u,s)]
\end{aligned}
\tag{5.12}
$$

where $p_\varphi(d^*|z,u,s)$ is the likelihood, $p_\alpha(z|d^*,u,s)$ is the posterior, $p_\beta(z|u,s)$ is the prior, and $D_{KL}$ is the Kullback-Leibler (KL) divergence. We use three neural networks to model these three probability distributions. An overview of our variational decipher is shown in Figure 5.3. We first use four recurrent neural networks to encode the (tweet, symbol, definition) pair in the dataset. Similar to what we do in the Seq2Seq model, there are two encoders for the hate symbol. One is at the word level and the other is at

Figure 5.3: The Variational Decipher. Note that this structure is used during training. During testing, the structure is slightly different. $d^*$ is the word embeddings of the definition. $x$ is the encoded definition. $c$ is the concatenation of the encoded tweet and hate symbol. $p$ and $p'$ are output distributions. $z$ is the latent variable. The definitions of other variables are the same as those in Figure 5.2.

the character level. The encoding of symbols and tweets are exactly the same as in our Seq2Seq model (see Eq. 5.2-5.4). The difference is that we also need to encode definitions for the Variational Decipher.

$$x, h_d = f_d(d^*) \tag{5.13}$$

Here, $f_d$ is the LSTM function. $x$ is the output of the LSTM at the last time step and $h_d$ is the hidden state of the LSTM at all time steps. The condition vector $c$ is the

concatenation of the encoded symbol words, symbol characters, and the tweet text:

$$c = c_u \oplus c_{sw} \oplus c_{sc} \tag{5.14}$$

We use multi-layer perceptrons (MLP) to model the posterior and the prior in the objective function. The posterior network and the prior network have the same structure and both output a probability distribution of latent variable $z$. The only difference is that the input of the posterior network is the concatenation of the encoded definition $x$ and the condition vector $c$ while the input of the prior network is only the condition vector $c$. Therefore, the output of the posterior network $p = p_\alpha(z|d^*, u, s)$ and the output of the prior network $p' = p_\beta(z|u, s)$. By assuming the latent variable $z$ has a multivariate Gaussian distribution, the actual outputs of the posterior and prior networks are the mean and variance: $(\mu, \Sigma)$ for the posterior network and $(\mu', \Sigma')$ for the prior network.

$$\mu, \Sigma = g(x \oplus c) \tag{5.15}$$

$$\mu', \Sigma' = g'(c) \tag{5.16}$$

$g$ is the MLP function of the posterior network and $g'$ is that of the prior network. During training, the latent variable $z$ is randomly sampled from the Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ and fed into the likelihood network. During testing, the posterior network is replaced by the prior network, so $z$ is sampled from $\mathcal{N}(\mu', \Sigma')$. The likelihood network is modeled by an RNN decoder with an attention mechanism, very similar to the decoder of our Seq2Seq model. The only difference lies in the input for the GRU. The decoder in our Variational Decipher model is to model the likelihood $p_\varphi(d^*|z, u, s)$, which is conditioned on the latent variable, tweet context, and the symbol. Therefore, for the Variational Decipher, the condition vector $c$ and the sampled latent variable $z$ are fed into the

decoder.

$$o_t, e_t = k(z \oplus c \oplus b_t, e_{t-1}) \tag{5.17}$$

$e_{t-1}$ is the hidden state of the RNN decoder at the last time step. $k$ is the GRU function. $o_t$ is its output and $e_t$ is its hidden state. Detailed decoding process and explanations are in Section 5.4.2.

---

**Algorithm 4** Train Variational Decipher

1: **function** TRAIN($U$)
2:      randomly initialize network parameters $\varphi$, $\alpha$, $\beta$;
3:      **for** $epoch = 1, E$ **do**
4:          **for** $(tweet, symbol, definition)$ in $U$ **do**
5:              get embeddings $u$, $s_w$, $s_c$, $d^*$;
6:              compute $x$, $c$ and $h$ with RNN encoders;
7:              compute $\mu$, $\Sigma$ with the posterior network;
8:              compute $\mu'$, $\Sigma'$ with the prior network;
9:              compute KL-divergence loss $L_{KL}$;
10:             sample $z = reparameterize(\mu, \Sigma)$;
11:             initialize the decoder state $e_0 = c$;
12:             $L_{REC} = 0$;
13:             **for** $t = 1, M$ **do**
14:                 compute attention weights $w_t$;
15:                 compute $o_t$, $e_t$ and $p(d_t|z, u, s)$;
16:                 $d_t = indmax(p(d_t|z, u, s))$;
17:                 $L_{REC} += BCE(d_t, d_t^*)$;
18:                 **if** $d_t$==EOS **then**
19:                     break;
20:                 **end if**
21:             **end for**
22:             update $\varphi$, $\alpha$, $\beta$ on $\mathcal{L} = L_{REC} + L_{KL}$;
23:         **end for**
24:     **end for**
25: **end function**

---

According to the objective function in Eq. 5.12, the loss function of the Variational

---
**Algorithm 5** Test Variational Decipher

---
1: **function** TEST($V$)
2:     **for** ($tweet, symbol, definition$) in $V$ **do**
3:         get embeddings $u$, $s_w$, $s_c$;
4:         compute $c$ and $h$ with RNN encoders;
5:         compute $\mu'$, $\Sigma'$ with the prior network;
6:         sample $z = reparameterize(\mu', \Sigma')$;
7:         initialize the decoder state $e_0 = c$;
8:         **for** $t = 1, M$ **do**
9:             compute attention weights $w$;
10:             compute $o_t$, $e_t$ and $p(d_t|z, u, s)$;
11:             $d_t = indmax(p(d_t|z, u, s))$;
12:             **if** $d_t$==EOS **then**
13:                 break;
14:             **end if**
15:         **end for**
16:     **end for**
17: **end function**

---

Decipher is as follows:

$$\mathcal{L} = L_{REC} + L_{KL}$$
$$= E_{z \sim p_\alpha(z|d^*, u, s)}[-\log p_\varphi(d^*|z, u, s)] + \quad (5.18)$$
$$D_{KL}[p_\alpha(z|d^*, u, s)||p_\beta(z|u, s)]$$

It consists of two parts. The first part $L_{REC}$ is called reconstruction loss. Optimizing $L_{REC}$ can push the sentences generated by the posterior network and the likelihood network closer to the given definitions. The second part $L_{KL}$ is the KL divergence loss. Optimizing this loss can push the output Gaussian Distributions of the prior network closer to that of the posterior network. This means we teach the prior network to learn the same knowledge learned by the posterior network, such that during testing time, when the referential definition $d^*$ is no longer available for generating the latent variable $z$, the prior network can still output a reasonable probability distribution over the latent

74

variable $z$. The complete training and testing process for the Variational Decipher is shown in Algorithm 4 and 5. $M$ is the predefined maximum length of the generated text. $BCE$ refers to the Binary Cross Entropy loss.

## 5.5 Experiments and Discussion

### 5.5.1 Experimental Settings

We use the dataset collected as described in Section 5.3 for training and testing. We randomly selected 2,440 tuples for testing and use the remaining 16,227 tuples for training. Note that there are no overlapping hate symbols between the training dataset $U$ and the testing dataset $D$.

We split the 2,440 tuples of the testing dataset $D$ into two separate parts, $D_s$ and $D_d$. $D_s$ consists of 1,681 examples and $D_d$ consists of 759 examples. In the first testing dataset $D_s$, although each hate symbol does not appear in the training dataset, the corresponding definition appears in the training dataset. In the second testing dataset $D_d$, neither the hate symbols nor the corresponding definitions appear in the training dataset. We do this split because deciphering hate symbols in these two cases has different levels of difficulty.

This split criterion means that for each hate symbol in $D_s$, there exists some symbol in the training dataset that has the same meaning but in different surface forms. For example, the hate symbol *wigwog* and *Wig Wog* have the same definition but one is in the training dataset, and the other is in the first testing dataset. We assume that such types of hate symbols share similar surface forms or similar tweet contexts. Therefore, the first testing dataset $D_s$ is to evaluate how well the model captures the semantic similarities among the tweet contexts in different examples or the similarities among different surface forms of a hate symbol.

Deciphering the hate symbols in the second testing dataset $D_d$ is more challenging. Both the unseen hate symbols and definitions require the model to have the ability to accurately capture the semantic information in the tweet context and then make a reasonable prediction. The second testing dataset $D_d$ is used to evaluate how well the model generalizes to completely new hate symbols.

For the Seq2Seq model, we use the negative log-likelihood loss for training. Both models are optimized using Adam optimizer [98]. The hyper-parameters of the two models are exactly the same. We set the maximum generation length $M = 50$. The hidden size of the encoders is 64. The size of the word embedding is 200 and that of character embedding is 100. The word embeddings and character embeddings are randomly initialized. Each model is trained for 50 epochs. We report the deciphering results of two models on three testing datasets $D$, $D_s$, and $D_d$.

## 5.5.2   Experimental Results

**Quantitative Results:** We use equally weighted BLEU score for up to 4-grams [99], ROUGE-L [100] and METEOR [101] to evaluate the decipherment results. The results are shown in Table 5.1. Figure 5.4 shows the BLEU score achieved by the two models on three testing datasets $D$, $D_s$, and $D_d$ during the training process. Both our Seq2Seq model and Variational Decipher achieve reasonable BLEU scores on the testing datasets. The Seq2Seq model outperforms the Variational Decipher on $D_s$ while Variational Decipher outperforms Seq2Seq on $D_d$. Note that $D_s$ is more than twice the size of $D_d$. Therefore, Seq2Seq outperforms Variational Decipher on the entire testing dataset $D$. The different performance of the two models on $D_s$ and $D_d$ is more obvious in Figure 5.4. The gap between the performance of the Seq2Seq model on $D_s$ and $D_d$ is much larger than that between the performance of the Variational Decipher on these two datasets.

| Dataset | Method | BLEU | ROUGE-L | METEOR |
|---------|--------|------|---------|--------|
| $D_s$ | Seq2Seq | **37.80** | **41.05** | **36.67** |
|  | VD | 34.77 | 32.96 | 31.03 |
| $D_d$ | Seq2Seq | 25.44 | 12.96 | **5.54** |
|  | VD | **28.38** | **14.01** | 5.41 |
| $D$ | Seq2Seq | **33.96** | **32.32** | **26.98** |
|  | VD | 32.75 | 27.00 | 23.16 |

Table 5.1: The BLEU, ROUGE-L, and METEOR scores on testing datasets. VD refers to the Variational Decipher. $D$ is the entire testing dataset. $D_s$ is the first part of $D$ and $D_d$ is the second part. The better results are in bold.



Figure 5.4: BLEU scores of two models on the testing dataset $D$, $D_s$, and $D_d$. The three dotted curves represent the performance of the Seq2Seq model while the three solid curves represent the performance of the Variational Decipher.

**Human Evaluation:** We employed crowd-sourced workers to evaluate the deciphering results of two models. We randomly sampled 100 items of deciphering results from $D_s$ and another 100 items from $D_d$. Each item composes a choice question and each choice question is assigned to five workers on Amazon Mechanical Turk. In each choice question, the workers are given the hate symbol, the referential definition, the original tweet, and two machine-generated plain texts from the Seq2Seq model and Variational Decipher. Workers are asked to select the more reasonable of the two results. In each choice question, the order of the results from the two models is permuted. Ties are permitted

| Dataset | Seq2Seq Lose | Seq2Seq Win | Tie |
|---------|-------------|-------------|-------|
| $D_s$ | 31.0% | 32.0% | 37.0% |
| $D_d$ | 30.5% | 22.0% | 47.5% |

Table 5.2: The results of human evaluation on two separate testing datasets $D_s$ and $D_d$.

for answers. We batch five items in one assignment and insert an artificial item with two identical outputs as a sanity check. The workers who fail to choose "tie" for that item are rejected from our test. The human evaluation results are shown in Table 5.2, which coincide with the results in Table 5.1 and Figure 5.4.

**Discussion:** When deciphering the hate symbols that have the same definitions as in the training dataset, the model can rely more on the surface forms of hate symbols than the tweet context to make a prediction because usually the hate symbols that share the same definitions also have similar surface forms. However, when it comes to the hate symbols with unseen definitions, simply relying on the surface forms cannot lead to a reasonable deciphering result. Instead, the model should learn the relationships between the context information and the definition of the symbol. Therefore, the different performances of two models on the two testing datasets $D_s$ and $D_d$ indicate that the Seq2Seq model is better at capturing the similarities among different surface forms of a hate symbol, while the Variational Decipher is better at capturing the semantic relationship between the tweet context and the hate symbol. The Sequence-to-Sequence model tries to capture such kinds of relationships by compressing all the context information into a fixed length vector, so its deciphering strategy is actually behavior cloning. On the other hand, the Variational Decipher captures such relationships by explicitly modeling the posterior and likelihood distributions. The modeled distributions provide higher-level semantic information compared to the compressed context, which allows the Variational Decipher to generalize better to the symbols with unseen definitions. This explains why the gap between the performance of the Seq2Seq model on two datasets is larger.

78

| Dataset | Symbol | tweet | Referential definition | Result of Seq2Seq | Result of VD |
|---|---|---|---|---|---|
| $D_s$ | *macaca* | What an ugly thing! And it's not because she's black, purple etc, it's because she's soooooooo ugly piece of *macaca* | Common French racist slur | Common French a racist | Common term black slur |
| $D_s$ | *closet homosexuals* | I wish the stupid Imams would just admit already that Muslim men are sex addicts & *closet homosexuals.* | A Homosexual who hasn't told anyone about his/her sexuality | A Homosexual who hasn't told anyone about and her sexuality | A Homosexual who hasn't told told anyone about his her sexuality |
| $D_s$ | *confederate flag* | done w/ this stupid sh*t. I support the *confederate flag* & will wear it. Will post pictures of it. Will fly it... | A flag that's usually flown in the south, most of the time flown to represent southern pride and heritage, but sometimes is flown to represent white power and racism | When s is or a any of but but | The from is and is or is or their south their to is or to |
| $D_d$ | *niggering* | More like call of n*gger: advanced *niggering* | The act of being a n*gger | A thing is to black to to to to | A thing but a n*gger n*gger n*gger |
| $D_d$ | *Heil Hitler* | *Heil Hitler*!!! RT @Am_Yuggio: 41.) Who is your role model ? :) | Its a gesture used by raising up the right hand with a straighten hand | A Schizophrenic of the the and Nazi Nazi by racist of Nazi more | Leader person and and Nazi and shit and |

Figure 5.5: Some example errors in the generated results of our Seq2Seq model and Variational Decipher.

## 5.5.3   Error Analysis

Figure 5.5 shows some example errors of the deciphering results of our Seq2Seq model and Variational Decipher. One problem with the deciphering results is that the generated sentences have a poor grammatical structure, as shown in Figure 5.5. This is mainly because the size of our dataset is small, and the models need a much larger corpus to learn the grammar. We anticipate that the generation performance will be improved with a larger dataset.

For the hate symbols in $D_s$, the deciphering results are of high quality when the length of referential definitions is relatively short. An example is *macaca*, a French slur shows in Figure 5.5. The deciphering result of the Seq2Seq model is close to the referential definition. As to the Variational Decipher, although the result is not literally the same as the definition, the meaning is close. *closet homosexuals* is another example. However, when the length of the referential definition increases, the performance of both models tends to be unsatisfactory, as the third example *confederate flag* shows. Although there exists the symbol *Confederate Flag* with the same definition in the training set, both

models fail on this example. One possible reason is that the complexity of generating the referential definition grows substantially with the increasing length, so when the tweet context and the symbol itself cannot provide enough information, the generation model cannot learn the relationship between the symbol and its definition.

Deciphering hate symbols in $D_d$ is much more challenging. Even for humans, deciphering completely new hate symbols is not a simple task. The two examples in Figure 5.5 show that the models have some ability to capture the semantic similarities. For the symbol *niggering*, the Variational Decipher generates the word *nigger* and Seq2Seq model generates *black*. For *Heil Hitler*, the Variational Decipher generates *leader person* and *Nazi*, while Seq2Seq also generates *Nazi*. Although these generated words are not in the definition, they still make some sense.

## 5.6   Conclusion

We propose a new task of learning to decipher hate symbols and create a symbol-rich tweet dataset. We split the testing dataset into two parts to analyze the characteristics of the Seq2Seq model and the Variational Decipher. The different performance of these two models indicates that the models can be applied to different scenarios of hate symbol deciphering. The Seq2Seq model outperforms the Variational Decipher for deciphering the hate symbols with similar definitions to that in the training dataset. This means the Seq2Seq model can better explain the hate symbols when Twitter users intentionally misspell or abbreviate common slur terms. On the other hand, the Variational Decipher tends to be better at deciphering hate symbols with unseen definitions, so it can be applied to explain newly created hate symbols on Twitter. Although both models show promising deciphering results, there still exists much room for improvement.

# Chapter 6

# Hate Speech Intervention

In the previous chapters, we primarily focus on developing NLP methods to automatically and effectively detect or analyze hate speech. In this chapter, we focus on the further action needed to calm and discourage individuals from using hate speech in the future. We propose a novel task of generative hate speech intervention, where the goal is to automatically generate responses to intervene during online conversations that contain hate speech. As a part of this work, we introduce two fully-labeled large-scale hate speech intervention datasets collected from Gab and Reddit[1]. These datasets provide conversation segments, hate speech labels, as well as intervention responses written by Mechanical Turk Workers. We also analyze the datasets to understand the common intervention strategies and explore the performance of common automatic response generation methods on these new datasets to provide a benchmark for future research.

---

[1]The datasets are publicly available here: `https://github.com/jing-qian/A-Benchmark-Dataset-for-Learning-to-Intervene-in-Online-Hate-Speech`

## 6.1   Introduction

To address the growing problem of online hate, an extensive body of work has focused on developing automatic hate speech detection models and datasets [5, 18, 10, 102, 103]. However, simply detecting and blocking hate speech or suspicious users often has limited ability to prevent these users from simply turning to other social media platforms to continue to engage in hate speech as can be seen in the large move of individuals blocked from Twitter to Gab [104]. What's more, such a strategy is often at odds with the concept of free speech. As reported by the Pew Research Center [105], "Despite this broad concern over online harassment, 45% of Americans say it is more important to let people speak their minds freely online; a slightly larger share (53%) feels that it is more important for people to feel welcome and safe online." The special rapporteurs representing the Office of the United Nations High Commissioner for Human Rights (OHCHR) have recommended that "The strategic response to hate speech is more speech." [106] They encourage to change what people think instead of merely changing what they do, so they advocate more speech that educates about cultural differences, diversity, and minorities as a better strategy to counter hate speech.

Therefore, in order to encourage strategies of countering online hate speech, we propose a novel task of generative hate speech intervention and introduce two new datasets for this task. Figure 6.1 illustrates the task. Our datasets consist of 5K conversations retrieved from Reddit and 12k conversations retrieved from Gab. Distinct from existing hate speech datasets, our datasets retain their conversational context and introduce human-written intervention responses. The conversational context and intervention responses are critical in order to build generative models to automatically mitigate the spread of these types of conversations.

To summarize, our contributions are three-fold:

Figure 6.1: An illustration of the hate speech conversation between User 1 and User 2 and the interventions collected for our datasets. The check and the cross icons on the right indicate a normal post and a hateful post. The utterance following the human icon is a human-written intervention, while the utterance following the computer icon is machine-generated.

- We introduce the generative hate speech intervention task and provide two fully-labeled hate speech datasets with human-written intervention responses.

- Our data is collected in the form of conversations, providing better context.

- The two data sources, Gab and Reddit, are not well studied for hate speech. Our datasets fill this gap.

Due to our data collecting strategy, all the posts in our datasets are manually labeled as hate or non-hate speech by Mechanical Turk workers, so they can also be used for the hate speech detection task. The performance of commonly-used classifiers on our datasets is shown in Section 6.6.

| | Source | #Posts | Conv. | Categories | Interv. |
|---|---|---|---|---|---|
| [18] | Twitter | 17k | No | racist, sexist, normal | No |
| [10] | Twitter | 25k | No | hateful, offensive, neither | No |
| [81] | Twitter | 35k | No | the worst, threats, hate speech, direct harassment, potentially offensive, non-harassment | No |
| [107] | Twitter | 9k | No | aggressive, bullying, spam, normal | No |
| [108] | Twitter, Reddit, The Gaurdian | 20k | No | harassment, non-harassment | No |
| [109] | Twitter | 100k | No | abusive, hateful, normal, spam | No |
| [5] | Yahoo! | 9k | No | anti-semitic, anti-black, anti-asian, anti-woman, anti-muslim, anti-immigrant, other-hate | No |
| [1] | Yahoo! | 2k | No | clean, hate, derogatory, profanity | No |
| [31] | Ask.fm | 85k | No | threat/blackmail, insult, defamation, sexual talk, curse/exclusion, defense, encouragement to the harasser | No |
| Ours | Reddit | 22k | Yes | hate, non-hate | Yes |
| Ours | Gab | 34k | Yes | hate, non-hate | Yes |

Table 6.1: Comparison of our datasets with previous hate speech datasets. Conv.: Conversation. Interv.: Intervention.

## 6.2    Related Work

In recent years, a few datasets for hate speech detection have been built and released by researchers. Most are collected from Twitter and are labeled using a combination of expert and non-expert hand labeling, or through machine learning assistance using a list of common negative words. It is widely accepted that labels can vary in their accuracy overall, though this can be mitigated by relying on a consensus rule to rectify disagreements in labels. A synopsis of these datasets can be found in Table 6.1.

In [18], the authors collect 17k tweets based on hate-related slurs and users. The tweets are manually annotated with three categories: sexist (20.0%), racist (11.7%), and normal (68.3%). Because the authors identified a number of prolific users during the initial manual search, the resulting dataset has a small number of users (1,236 users)

involved, causing a potential selection bias. This problem is most prevalent on the 1,972 racist tweets, which are sent by only 9 Twitter users. To avoid this problem, we did not identify suspicious user accounts or utilize user information when collecting our data.

A similar strategy is used in [10], which combines the utilization of hate keywords and suspicious user accounts to build a dataset from Twitter. But different from [18], this dataset consists of 25k tweets randomly sampled from the 85.4 million posts of a large number of users (33,458 users). This dataset is proposed mainly to distinguish hateful and offensive language, which tend to be conflated by many studies.

In [81], the authors focus on online harassment on Twitter and propose a fine-grained labeled dataset with 6 categories. A large Twitter dataset with 100k tweets is introduced in [109]. Despite the large size of this dataset, the ratio of hateful tweets is relatively low (5%). Thus the size of the hateful tweets is around 5k in this dataset, which is not significantly larger than that of the previous datasets.

The dataset introduced in [107] is different from the other datasets as it investigates the behavior of hate-related users on Twitter, instead of evaluating hate-related tweets. The large majority of the 1.5k users are labeled as spammers (31.8%) or normal (60.3%). Only a small fraction of the users are labeled as bullies (4.5%) or aggressors (3.4%).

While most datasets are from single sources, a dataset with a combination of Twitter (58.9%), Reddit, and The Guardian is introduced in [108]. In total 20,432 unique comments were obtained with 4,136 labeled as harassment (20.2%) and 16,296 as non-harassment (79.8%).

Since most of the publicly available hate speech datasets are collected from Twitter, previous research on hate speech mainly focuses on Twitter posts or users [18, 97, 29, 19, 10]. While there are several studies on the other sources, such as Instagram [6], Yahoo! [5, 1], and Ask.fm [31], the hate speech on Reddit and Gab is not widely studied. What's more, all the previous hate speech datasets are built for the classification or detection of

hate speech from a single post or user on social media, ignoring the context of the post and intervention methods needed to effectively calm down the users and diffuse negative online conversations.

## 6.3    Data

### 6.3.1    Ethics

Our study got approval from our Internal Review Board. Workers were warned about the offensive content before they read the data and they were informed by our instructions to feel free to quit the task at any time if they are uncomfortable with the content. Additionally, all personally identifiable information such as user names is masked in the datasets.

### 6.3.2    Data Filtering

**Reddit:** To retrieve high-quality conversational data that would likely include hate speech, we referenced the list of the whiniest most low-key toxic subreddits [110]. Skipping the three subreddits that have been removed, we collect data from ten subreddits: *r/DankMemes, r/Imgoingtohellforthis, r/KotakuInAction, r/MensRights, r/MetaCanada, r/MGTOW, r/PussyPass, r/PussyPassDenied, r/The_Donald*, and *r/TumblrInAction*. For each of these subreddits, we retrieve the top 200 hottest submissions using Reddit's API. To further focus on conversations with hate speech in each submission, we use hate keywords [103] to identify potentially hateful comments and then reconstructed the conversational context of each comment. This context consists of all comments preceding and following a potentially hateful comment. Thus for each potentially hateful comment, we rebuild the conversation where the comment appears. Figure 6.2 shows an example

| Conversation | Hate Speech | Human-Written Intervention Responses |
|---|---|---|
| 1. User 1: United Kingdom: 'Schoolboy, 15, given detention for backing UKIP during classroom debate'<br>2.    User 2: The education system is full of re\*\*\*ds! Yes, most school teachers are ret\*\*\*ed lefties! Teach your children to laugh at these ret\*\*\*ed lefties!<br>3.       User 3: Asking a teacher to not be a leftist is like asking a medieval munk to question the Pope.<br>4.       User 4: The Jews are like Sjws, they infest everything. | 2, 4 | ➢ Use of this language is not tolerated and it is uncalled for.<br>➢ Use of the slurs and insults here is unacceptable in our discourse as it demeans and insults and alienates others.<br>➢ I recommend that you research the holocaust, you might change your opinion. |

Figure 6.2: An example of the aggregated data. The first column is the conversation text. Indexes are added to each post. Indentations before each post indicate the structure of replies. The second column is the indexes of the human-labeled hateful post. Each bullet point in the third column is a human-written response.

of the collected conversation, where the second comment contains a hate keyword and is considered potentially hateful. Because a conversation may contain more than one comments with hate keywords, we removed any duplicated conversations.

**Gab:** We collect data from all the Gab posts in October 2018. Similar to Reddit, we use hate keywords [103] to identify potentially hateful posts, rebuild the conversation context, and clean duplicate conversations.

## 6.3.3   Crowd-Sourcing

After we collected the conversations from both Reddit and Gab, we presented this data to Mechanical Turk workers to label and create intervention suggestions. In order not to over-burden the workers, we filtered out conversations consisting of more than 20 comments. Each assignment consists of 5 conversations. For Reddit, we also present the title and content of the corresponding submission in order to give workers more information about the topic and context. For each conversation, a worker is asked to answer two questions:

- Q1: Which posts or comments in this conversation are hate speech?

- Q2: If there exists hate speech in the conversation, how would you respond to

87

intervene? Write down a response that can probably hold it back (word limit: 140 characters).

If the worker thinks no hate speech exists in the conversation, then the answers to both questions are "n/a". To provide context, the definition of hate speech from Facebook [111]: *"We define hate speech as a direct attack on people based on what we call protected characteristics — race, ethnicity, national origin, religious affiliation, sexual orientation, caste, sex, gender, gender identity, and serious disease or disability."* is presented to the workers. Also, to prevent workers from using hate speech in the response or writing responses that are too general, such as *"Please do not say that"*, we provide additional instructions and rejected examples.

## 6.3.4  Data Quality

Each conversation is assigned to three different workers. To ensure data quality, we restrict the workers to be in an English-speaking countries including Australia, Canada, Ireland, New Zealand, the United Kingdom, and the United States, with a HIT approval rate higher than 95%. Excluding the rejected answers, the collected data involves 926 different workers. The final hate speech labels (answers to Q1) are aggregated according to the majority of the workers' answers. A comment is considered hate speech only when at least two out of the three workers label it as hate speech. The responses (answers to Q2) are aggregated according to the aggregated result of Q1. If the worker's label to Q1 agrees with the aggregated result, then their answer to Q2 is included as a candidate response to the corresponding conversation but is otherwise disregarded. See Figure 6.2 for an example of the aggregated data.

## 6.4   Data Analysis

### 6.4.1   Statistics

From Reddit, we collected 5,020 conversations, including 22,324 comments. On average, each conversation consists of 4.45 comments and the length of each comment is 58.0 tokens. 5,257 of the comments are labeled as hate speech and 17,067 are labeled as non-hate speech. A majority of the conversations, 3,847 (76.6%), contain hate speech. Each conversation with hate speech has 2.66 responses on average, for a total of 10,243 intervention responses. The average length of the intervention responses is 17.96 tokens.

From Gab, we collected 11,825 conversations, consisting of 33,776 posts. On average, each conversation consists of 2.86 posts and the average length of each post is 35.6 tokens. 14,614 of the posts are labeled as hate speech and 19,162 are labeled as non-hate speech. Nearly all the conversations, 11,169 (94.5%), contain hate speech. 31,487 intervention responses were originally collected for conversations with hate speech, or 2.82 responses per conversation on average. The average length of the intervention responses is 17.27 tokens.

Compared with the Gab dataset, there are fewer conversations and comments in the Reddit dataset, comments and conversations are longer, and the distribution of hate and non-hate speech labels is more imbalanced. Figure 6.3 illustrates the distributions of the top 10 keywords in the hate speech collected from Reddit and Gab separately. The Gab dataset and the Reddit dataset have similar popular hate keywords, but the distributions are very different. All the statistics shown above indicate that the characteristics of the data collected from these two sources are very different, thus the challenges of doing detection or generative intervention tasks on the dataset from these sources will also be different.

Figure 6.3: The distributions of the top 10 keywords in the hate speech collected from Reddit and Gab. Hate keywords are masked.

## 6.4.2  Intervention Strategies

Removing duplicates, there are 21,747 unique intervention responses in the aggregated Gab dataset and 7,641 in the aggregated Reddit dataset. Despite the large diversity of the collected responses for intervention, we find workers tend to have certain strategies for intervention.

**Identify Hate Keywords**: One of the most common strategies is to identify the inappropriate terms in the post and then urge the user to stop using that work. For example, *"The C word and language attacking gender is unacceptable. Please refrain from future use."* This strategy is often used when the hatred in the post is mainly conveyed by specific hate keywords.

**Categorize Hate Speech:** This is another common strategy used by the workers. The workers classify hate speech into different categories, such as racist, sexist, homophobic, etc. This strategy is often combined with identifying hate keywords or targets of hatred. For example, *"The term ""fa\*\*ot"" comprises homophobic hate, and as such is not*

*permitted here."*

**Positive Tone Followed by Transitions:** This is a strategy where the response consists of two parts combined with a transitional word, such as "but" and "even though". The first part starts with affirmative terms, such as "I understand", "You have the right to", and "You are free to express", showing kindness and understanding, while the second part is to alert the users that their post is inappropriate. For example, *"I understand your frustration, but the term you have used is offensive towards the disabled community. Please be more aware of your words."*. Intuitively, compared with the response that directly warns, this strategy is likely more acceptable for the users and is more likely to calm down a quarrel full of hate speech.

**Suggest Proper Actions:** Besides warning and discouraging the users from continuing hate speech, workers also suggest the actions that the user should take. This strategy can either be combined with other strategies mentioned above or be used alone. In the latter case, a negative tone can be greatly alleviated. For example, *"I think that you should do more research on how resources are allocated in this country."*

## 6.5  Generative Intervention

Our datasets can be used for various hate speech tasks. In this chapter, we focus on generative hate speech intervention. The goal of this task is to generate a response to hate speech that can mitigate its use during a conversation. The objective can be formulated as the following equation:

$$Obj = \max \sum_{(c,r)\in D} \log p(r|c) \tag{6.1}$$

91

where $c$ is the conversation, $r$ is the corresponding intervention response, and $D$ is the dataset. This task is closely related to response generation and dialog generation, though several differences exist including dialog length, language cadence, and word imbalances. As a baseline, we chose the most common methods of these two tasks, such as Seq2Seq and VAE, to determine the initial feasibility of automatically generating intervention responses. More recent Reinforcement Learning methods for dialog generation [112] can also be applied to this task with slight modification. Future work will explore more complex, and unique models.

Similar to [112], a generative model is considered as an agent. However, different from dialog generation, generative intervention does not have multiple turns of utterance, so the action of the agent is to select a token in the response. The state of the agent is given by the input posts and the previously generated tokens. Another result due to this difference is that the rewards with regard to ease of answering or information flow do not apply to this case, but the reward for semantic coherence does. Therefore, the reward of the agent is:

$$rw(c, r) = \lambda_1 \log p(r|c) + \lambda_2 \log p_{back}(c|r) \tag{6.2}$$

where $rw(c, r)$ is the reward with regard to the conversation $c$ and its reference response $r$ in the dataset. $p(r|c)$ denotes the probability of generating response $r$ given the conversation $c$, and $p_{back}(c|r)$ denotes the backward probability of generating the conversation based on the response, which is parameterized by another generation network. The reward is a weighted combination of these two parts, which are observed after the agent finishes generating the response. For details, please refer to the original work [112].

## 6.6 Experiments and Discussion

We evaluate the commonly-used detection and generation methods with our dataset. Due to the different characteristics of the data collected from the two sources (Section 6.4), we treat them as two independent datasets.

### 6.6.1 Experimental Settings

For **binary hate speech detection**, we experimented with the following four different methods.

**Logistic Regression (LR):** We evaluate the Logistic Regression model with L2 regularization. The penalty parameter C is set to 1. The input features are the Term Frequency Inverse Document Frequency (TF-IDF) values of up to 2-grams.

**Support Vector Machine (SVM):** We evaluate the SVM model with linear kernels. We use L2 regularization and the coefficient is 1. The features are the same as in LR.

**Convolutional Neural Network (CNN):** We use the CNN model for sentence classification proposed by [41] with default hyperparameters. The word embeddings are randomly initialized (CNN in Table 6.2) or initialized with pretrained Word2Vec [113] embeddings on Google News (CNN* in Table 6.2).

**Recurrent Neural Network (RNN):** The model we evaluated consists of a 2-layer bidirectional Gated Recurrent Unit (GRU) [83] followed by a linear layer. Same as for CNN, we report the performance of RNN with two different settings of the word embeddings.

The methods are evaluated on testing data randomly selected from the dataset with a ratio of 20%. The input data is not manipulated to manually balance the classes for any of the above methods. Therefore, the training and testing data retain the same distribution as the collected results (Section 6.4). The methods are evaluated using the F-1 score,

Precision-Recall (PR) AUC, and Receiver-Operating-Characteristic (ROC) AUC.

For **generative hate speech intervention**, we evaluated the following three methods.

**Seq2Seq** [84, 83]: The encoder consists of 2 bidirectional GRU layers. The decoder consists of 2 GRU layers followed by a 3-layer MLP (Multi-Layer Perceptron).

**Variational Auto-Encoder (VAE)** [54]: The structure of the VAE model is similar to that of the Seq2Seq model, except that it has two independent linear layers followed by the encoder to calculate the mean and variance of the distribution of the latent variable separately. We assume the latent variable follows a multivariate Gaussian Distribution. KL annealing [114] is applied during training.

**Reinforcement Learning (RL):** We also implement the Reinforcement Learning method described in Section 6.5. The backbone of this model is the Seq2Seq model, which follows the same Seq2Seq network structure described above. This network is used to parameterize the probability of a response given the conversation. Besides this backbone Seq2Seq model, another Seq2Seq model is used to generate the backward probability. This network is trained in a similar way as the backbone Seq2Seq model, but with a response as input and the corresponding conversation as the target. In our implementation, the function of the first part of the reward $(\log p(r|c))$ is conveyed by the MLE loss. A curriculum learning strategy is adopted for the reward of $\log p_{back}(c|r)$ as in [115]. Same as in [112] and [115], a baseline strategy is employed to estimate the average reward. We parameterize it as a 3-layer MLP.

The Seq2Seq model and VAE model are evaluated under two different settings. In one setting, the input for the generative model is the complete conversation, while in the other setting, the input is the filtered conversation, which only includes the posts labeled as hate speech. The filtered conversation was necessary to test the Reinforcement Learning model, as it is too challenging for the backward model to reconstruct the complete

conversation based only on the intervention response.

In our experiments on the generative hate speech intervention task, we do not consider conversations without hate speech. The testing dataset is then randomly selected from the resulting dataset with a ratio of 20%. Since each conversation can have multiple reference responses, we dis-aggregate the responses and construct a pair (conversation, reference response) for each of the corresponding references during training. Teacher forcing is used for each of the three methods. The automatic evaluation metrics include BLEU [99], ROUGE-L [100], and METEOR [101].

In order to validate and compare the quality of the generated results from each model, we also conducted human evaluations as previous research has shown that automatic evaluation metrics often do not correlate with human preference [116]. We randomly sampled 450 conversations from the testing dataset. We then generated responses using each of the above models trained with the filtered conversation setting. In each assignment, a Mechanical Turk worker is presented with 10 conversations, along with corresponding responses generated by the three models. For each conversation, the worker is asked to evaluate the effectiveness of the generated intervention by selecting a response that can best mitigate hate speech. 9 of the 10 questions are filled with the sampled testing data and the generated results, while the other is artificially constructed to monitor response quality. After selecting the 10 best mitigation measures, the worker is asked to select which of the three methods has the best diversity of responses over all the 10 conversations. Ties are permitted for answers. Assignments that failed the quality check are rejected.

| Conversation | Seq2Seq | VAE | RL |
|---|---|---|---|
| User 1: Ok. Keep White Knighting for a sand ni**er who f**ks tons of White women. You're entitled to your opinion. | This is is not allowed and un time to treat people by their skin color | Use of the r-word is unacceptable in our discourse as it demeans and insults people with mental disabilities | I do not use derogatory racial slurs |
| User 1: Feminists are f**king ret***ed, why are they so pissed about manspreading shit makes no difference in space. I'm gonna start throwing things in womens cleavage as they walk by. | Using talking about someones character to describe a woman is insulting to hateful | I don't think using words that are sexist in nature contribute to a productive conversation | Your post uses speech that denigrates sex and gender, which is in violation of our Community Guidelines and TOS Agreement |
| User 1: Only a nig****ick could reach that c**t<br>    User 2: I bet a midget football team could hide under her overhang during a rainstorm and not get wet.<br>    User 3: Ni**er probably uses her as a waterbed | Use of the c-word is unacceptable in our discourse as it demeans and insults women | Please do not use derogatory language for intellectual disabilities | If you don't agree with you, there 's no need to resort to name calling |

Figure 6.4: Examples of the generated intervention responses. The hateful terms in the conversation are masked.

| Dataset | Gab | | | Reddit | | |
|---|---|---|---|---|---|---|
| Metric | F1 | PR | ROC | F1 | PR | ROC |
| LR | 88.2 | 94.5 | 95.4 | 64.7 | 80.4 | 91.4 |
| SVM | 88.6 | 94.7 | 95.6 | 75.7 | **81.1** | **92.0** |
| CNN | 87.5 | 92.8 | 92.6 | 74.8 | 76.8 | 87.5 |
| RNN | 87.6 | 93.9 | 94.2 | 71.7 | 76.1 | 88.6 |
| CNN* | **89.6** | **95.2** | **95.8** | 76.9 | 80.1 | 90.9 |
| RNN* | 89.3 | 94.8 | 95.5 | **77.5** | 79.4 | 90.6 |

Table 6.2: Experimental results for the detection task. PR is Precision-Recall AUC and ROC is ROC AUC. The models marked with * use pretrained Word2Vec embeddings. The best results are in bold.

## 6.6.2 Experimental Results and Discussion

The experimental results of the detection task and the generative intervention task are shown in Table 6.2 and Table 6.3 separately. The results of the human evaluation are shown in Table 6.4. Figure 6.4 shows examples of the generated responses.

As shown in Table 6.2 and 6.3, all the classification and generative models perform better on the Gab dataset than on the Reddit dataset. We think this stems from the datasets' characteristics. First, the Gab dataset is larger and has a more balanced category distribution than the Reddit dataset. Therefore, it is inherently more challenging to train a classifier on the Reddit dataset. Further, the average lengths of the Reddit posts

| Dataset | Gab | | | | | | Reddit | | | | | |
|---------|-----|--|--|--|--|--|--------|--|--|--|--|--|
| Inp. set. | Complete | | | Filtered | | | Complete | | | Filtered | | |
| Metric | B | R | M | B | R | M | B | R | M | B | R | M |
| Seq2Seq | **13.2** | **33.8** | 23.0 | **15.0** | **34.2** | 23.6 | 5.5 | **29.5** | 19.5 | 5.9 | 28.2 | 20.0 |
| VAE | 12.2 | 32.5 | **23.4** | 12.4 | 32.8 | 21.8 | **6.8** | 29.0 | **20.2** | **7.0** | **29.1** | **20.1** |
| RL | - | - | - | 14.5 | 33.1 | **23.9** | - | - | - | 4.4 | **29.1** | 18.7 |

Table 6.3: Experimental results for generative intervention task. Inp. set.: Input Setting (Section 6.6.1). B: BLEU. R: ROUGE-L. M: METEOR. The best results are in bold.

| Dataset | Gab | | Reddit | |
|---------|-----|--|--------|--|
| Metric | Eff. | Div. | Eff. | Div. |
| Seq2Seq Wins | 22.4 | 28.0 | **31.1** | **34.0** |
| VAE Wins | 20.0 | 6.0 | 26.0 | 4.0 |
| RL Wins | **41.6** | **40.0** | 30.0 | 30.0 |
| Tie | 16.0 | 26.0 | 12.9 | 32.0 |

Table 6.4: Human evaluation results. Table values are the percentage of the answers. Eff.: Effectiveness, evaluates how well the generated responses can mitigate hate speech. Div: Diversity, evaluates how many different responses are generated. The best results are in bold.

and conversations are much larger than those of Gab, potentially making the Reddit input nosier than the Gab input for both tasks. On both the Gab and Reddit datasets, the SVM classifier and the LR classifier achieved better performance than the CNN and RNN models with randomly initialized word embeddings. A possible reason is that without pretrained word embeddings, the neural network models tend to overfit on the dataset.

For the generative intervention task, three models perform similarly on all three automatic evaluation metrics. As expected, the Seq2Seq model achieves higher scores with filtered conversation as input. However, this is not the case for the VAE model. This indicates that the two models may have different capabilities to capture important information in conversations.

As shown in Table 6.3, applying Reinforcement Learning does not lead to higher scores on the three automatic metrics. However, human evaluation (Table 6.4) shows that the

RL model creates responses that are potentially better at mitigating hate speech and are more diverse, which is consistent with [112]. There is a larger performance difference with the Gab dataset, while the effectiveness and the diversity of the responses generated by the Seq2Seq model and the RL model are quite similar on the Reddit dataset. One possible reason is that the size of the training data from Reddit (around 8k) is only 30% the size of the training data from Gab. The inconsistency between the human evaluation results and the automatic ones indicates the automatic evaluation metrics listed in Table 6.3 can hardly reflect the quality of the generated responses. As mentioned in Section 6.4, annotators tend to have strategies for intervention. Therefore, generating the common parts of the most popular strategies for all the testing inputs can lead to high scores on these automatic evaluation metrics. For example, generating *"Please do not use derogatory language."* for all the testing Gab data can achieve 4.2 on BLEU, 20.4 on ROUGE, and 18.2 on METEOR. However, this response is not considered high-quality because it is almost a universal response to all hate speech, regardless of the context and topic.

Surprisingly, the responses generated by the VAE model have much worse diversity than the other two methods according to human evaluation. As indicated in Figure 6.4, the responses generated by VAE tend to repeat the responses related to some popular hate keyword. For example, *"Use of the r-word is unacceptable in our discourse as it demeans and insults people with mental disabilities."* and *"Please do not use derogatory language for intellectual disabilities."* are the generated responses for a large part of the Gab testing data. According to Figure 6.3, insults toward disabilities are the largest portion of the dataset, so we suspect that the performance of the VAE model is affected by the imbalanced keyword distribution.

The sampled results in Figure 6.4 show that the Seq2Seq and the RL model can generate reasonable responses for intervention. However, as is to be expected with

machine-generated text, in the other human evaluation we conducted, where Mechanical Turk workers were also presented with sampled human-written responses alongside the machine-generated responses, the human-written responses were chosen as the most effective and diverse option a majority of the time (70% or more) for both datasets. This indicates that there is significant room for improvement while generating automated intervention responses.

In our experiments, we only utilized the text of the posts, but more information is available and can be utilized, such as the user information and the title of a Reddit submission.

## 6.7    Conclusion

Towards the end goal of mitigating the problem of online hate speech, we propose the task of generative hate speech intervention and introduce two fully-labeled datasets collected from Reddit and Gab, with crowd-sourced intervention responses. The performance of the three generative models: Seq2Seq, VAE, and RL, suggests ample opportunity for improvement. We make our datasets freely available to facilitate further exploration of hate speech intervention and better models for generative intervention.

# Part II

# Language Model Detoxification

# Chapter 7

# Language Model Detoxification in Freeform Text Generation

In the previous chapters, we have talked about the automatic post-processing of hate speech. In this part of the dissertation, we tackle the toxic degeneration problem of large pretrained Language Models (LMs) by controlling the generation of LMs. Previous work in this field has focused on directly fine-tuning the language model or utilizing an attribute discriminator. In this chapter, we propose a novel lightweight framework for controllable GPT2 [14] generation, which utilizes a set of small attribute-specific vectors, called prefixes [117], to steer natural language generation. Different from [117], where each prefix is trained independently, we take the relationship among prefixes into consideration and train multiple prefixes simultaneously. We propose a novel supervised method and also an unsupervised method to train the prefixes for detoxification. Experimental results show that our methods can guide generation towards low toxicity while keeping high linguistic quality.

## 7.1   Introduction

Currently, the prevalent LMs are pretrained on large amounts of text from the web. Large-scale pretraining enables these LMs to achieve a good estimation of the next token probability, resulting in coherent texts. However, relying on web text for pretraining also results in degeneration and biased behavior. Previous research [7] shows that LMs pretrained on problematic web content can easily degenerate into toxicity, even without explicitly toxic prompts. One direction to alleviate this problem is to apply post-processing. Our work on automatic toxic speech detection (Chapter 2) can be applied to filter out the generations with toxic speech, and the LM repeatedly does the generation until the classifier finds a safe one. However, this strategy is inefficient and there may exist cases where there is no safe choice within a fixed number of generations [24].

In order to circumvent this limitation, recent research has focused on controlling the generation of these Natural Language Generation (NLG) models from the source. The goal of controllable NLG is to guide generation towards the desired attributes in the concerned aspects of the text, such as toxicity, sentiment, or topic. LM Detoxification can be considered as an application of the control method to infuse the safe language attributes. Previous work has focused on directly fine-tuning the existing models [118, 119, 120] or using a discriminator to guide generation [25, 26, 121]. CTRL [118] achieves controllability at the expense of training a large conditional LM. GeDi [26] also trains conditional LMs but uses them as discriminators to guide generation, introducing additional 345M parameters. PPLM [25] guides generation by iteratively updating the LM's hidden activations. However, this decoding strategy is extremely computationally intensive, resulting in a slow generation speed [7].

Prefix-tuning [117] proposes to optimize a prefix, which is a small continuous task-specific vector, as a lightweight alternative to fine-tuning an NLG task, such as table-

Figure 7.1: A comparison of prefix-tuning [117] (top) and our framework (bottom) on detoxification. The solid arrows show the training process, while the dashed ones show the inference (generation) process. In our proposed framework, the training can be supervised or unsupervised.

to-text generation or summarization. Inspired by prefix-tuning [117], we propose to use prefixes, a set of small continuous attribute-specific vectors, to steer NLG. Compared with using an attribute model or a generative discriminator [25, 26], using learned prefixes to achieve controllability has the following benefits. First, it introduces fewer additional parameters ($\sim$0.2%-2% of GPT2 parameters in our experiments). Second, using prefixes keeps the inference speed comparable to that of the original GPT2 model.

In a general sense, prefix-tuning [117] can be considered as controlling the generation of language models. Prefix-tuning views each prefix as an independent control task thus trains each prefix separately (top in Figure 7.1). However, in the web text, non-toxic speech and toxic speech are very imbalanced in both diversity and quantity. The training examples of non-toxic speech may not teach the model about what is considered toxic and what to avoid during generation. We think that having the model learn about what is considered toxic is also important for detoxification. Therefore, we propose a novel supervised method and a novel unsupervised one in our framework, which takes

the relationship among prefixes into consideration and trains two prefixes simultaneously with novel training objectives, as illustrated in Figure 7.1.

Experimental results on detoxification show that our proposed methods can guide generation towards low toxicity while keeping high linguistic quality, even when only several dozen labeled examples are available.

Our main contributions are as follows:

- We propose a novel framework that utilizes prefixes with frozen LMs as a lightweight alternative for controllable GPT2 generation.

- We propose a supervised method and an unsupervised method with novel objectives for prefix training, where the relationship among prefixes is considered and multiple prefixes are trained simultaneously.

- Experimental results show that our methods can effectively guide generation towards low toxicity while keeping high linguistic quality.

## 7.2   Related Work

With the rapid development of NLG, researchers have attempted to interpolate the controlling factors into generation. In [120], the stylistic aspects of the generated text are controlled with a conditioned RNN (Recurrent Neural Network) LM. In [121], the authors compose a committee of discriminators to guide an RNN generator towards the generations with the desired linguistic quality. In [119], the authors aim at controlling the sentiment and tense of the generated text by combining variational auto-encoders (VAE) and attribute discriminators.

More recently, with the advent of Transformers and large pretrained language models, such as GPT2, an extensive body of work has focused on controlling the generation

of these Transformer-based models. CTRL [118] trains a 1.63 billion-parameter conditional transformer LM from scratch with 55 attribute control codes to guide generation. However, this method is expensive and lacks flexibility since the control codes are fixed. PPLM [25] addresses these limitations by developing a plug-and-play model which leverages an attribute discriminator to perturb the LM's hidden activations. However, updating gradients at the token level results in slow inference. Instead of updating the hidden activations, later work [26, 122, 123] introduces generative discriminators to re-weight the next token distributions on the fly during inference, thus improving the inference speed. Our work is mostly related to [124, 117]. In [124], the authors use a pretrained LM followed by an attribute alignment function to encode the tokens of the target attributes and the resulting hidden states are used to control generation. Different from their work, we do not take the tokens of the target attributes as input. Instead, we directly train a set of parameters, which acts as the prepended hidden states of GPT2, to control generation. Avoiding using attribute tokens can circumvent the problems when it is difficult to describe the desired attribute with only one word. Besides, in [124], the authors focus on attribute disentanglement, which is not a focus in our work, so our training methods are different. Prefix-tuning [117] can, in a general sense, be viewed as controlling the generation of LMs, where the LM is controlled to depict a specific NLG task, while in this work, the LM is controlled to carry a specific attribute in a generation. Besides, our proposed methods for prefix training are different from prefix-tuning, as stated in Section 7.1.

## 7.3   Methods

Our method uses prefixes to guide GPT2 generation, where a prefix is a continuous attribute-specific vector prepended to the activations of GPT2. Prefixes are free pa-

Figure 7.2: An illustration of the GPT2 generation process unfolded through time, controlled by a negative toxicity prefix $H_0 = H_\theta[0, :, :]$. *"You are"* is the given prompt. *"a nice"* is the generated completion.

rameters denoted as $H_\theta$. Different from prefix-tuning [117], where each prefix is trained independently, we consider the relationship among attributes and train multiple prefixes simultaneously, so $H_\theta$ is of dimension $N \times M \times D$, where $N$ is the number of prefixes. In our detoxification task, $N$ equals 2. $M$ is the length of a prefix. $D = 2 \times L \times E$ is the dimension of the activation in GPT2, where $L$ is the number of transformer layers, $E$ is the hidden size, and 2 indicates one key vector and one value vector. Following prefix-tuning, we reparametrize $H_\theta[i, j, :] = W_i H'_\theta[i, j, :]$ by a smaller parameter ($H'_\theta$) composed with a large matrix ($W_i$). After the training finishes, only $H_\theta$ needs to be saved for generation while $W$ and $H'_\theta$ can be discarded. Since the GPT2 parameters are kept frozen during training, they do not need to be saved either. Figure 7.2 shows an example of the generation process under the control of a trained prefix. The prefixes can be trained in a supervised or unsupervised way.

## 7.3.1   Supervised Method

Suppose the concerned aspect has the attribute set $Y$, each training example is a pair of $(x, y)$ where $x$ is the input text and $y \in Y$ is the attribute label of $x$. Note that the attribute label also indicates the ground truth index of the prefix in $H_\theta$, so $y$ also refers to the prefix index in the following description. As mentioned in Section 7.1, we introduce an additional discriminative loss to train multiple prefixes simultaneously. Therefore, the training loss $\mathcal{L}_{sup}$ is a weighted sum of the language model loss $\mathcal{L}_{LM}$ and the discriminative loss $\mathcal{L}_d$:

$$\mathcal{L}_{sup} = \omega_1 \mathcal{L}_{LM} + \omega_2 \mathcal{L}_d \tag{7.1}$$

$$\mathcal{L}_{LM} = -\sum_{t=1}^{T} \log p(x_t | x_{<t}, y) \tag{7.2}$$

$$\mathcal{L}_d = -\log \frac{p(y)p(x|y)}{\sum_{y' \in Y} p(y')p(x|y')} \tag{7.3}$$

The computation of $\log p(x_t | x_{<t}, y)$ is parameterized as $\log p_{\theta, \gamma}(x_t | x_{<t}, H_\theta[y, :, :])$, where $\gamma$ is the set of fixed GPT2 parameters, and $\theta$ represents learnable prefix parameters. $\log p(x|y) = \sum_t \log p(x_t | x_{<t}, y)$, so the parameterization of $\log p(x|y)$ is the sum of $\log p_{\theta, \gamma}(x_t | x_{<t}, H_\theta[y, :, :])$ over $t$.

Note that each prefix can be trained independently using $\mathcal{L}_{LM}$ alone, which would be the same as prefix-tuning [117]. Intuitively, prefixes trained by $\mathcal{L}_{LM}$ are infused with the information of what is encouraged to generate. However, we observe that in controllable NLG, especially in detoxification, it is helpful to also infuse a prefix with the information of what is discouraged to generate. Given a training example $(x, y)$, the prefix $H_\theta[y, :, :]$ should be optimized towards generating $x$, while the other prefixes should be discouraged to generate $x$. To achieve this goal, all the prefixes in $H_\theta$ should be trained simultaneously. Therefore, the discriminative loss $\mathcal{L}_d$ is introduced. As in Equation 7.3,

Figure 7.3: An illustration of the supervised training method. $H_0$ is the prefix of negative toxicity. $H_1$ is the prefix of positive toxicity. Note that training without $\mathcal{L}_d$ is equivalent to prefix-tuning [117], where $H_0$ and $H_1$ are trained separately. $p_i$ and $p_i'$ are the next token probabilities. $p_i = p(x_i|x_{<i}, H_1)$ and $p_i' = p(x_i|x_{<i}, H_0)$. The GPT2 is pretrained, and its parameters are frozen.



Figure 7.4: An illustration of the unsupervised training method. $H_\theta$ denotes the 2 prefixes. $z$ is the latent variable indicating the index of the prefix corresponding to the input text $x$. $\bar{z}$ is the latent variable indicating the index of the opposite prefix. $\otimes$ is matrix multiplication. $\mathcal{L}_{KL}$ is not shown in this figure for clarity.

optimizing $\mathcal{L}_d$ improves the attribute alignment $p(y|x)$ by increasing $p(x|y)$ and lowering $p(x|\bar{y})$, $\bar{y} \in Y \backslash \{y\}$ at the same time. We assume uniform prior, so $p(y)$ and $p(y')$ can be canceled out in Equation 7.3. Figure 7.3 illustrates the training process with two prefixes.

## 7.3.2   Unsupervised Method

In the unsupervised setting, we assume the attribute set $Y$ of the concerned aspect is known. The training example consists of input text $x$ only. The attribute label $y$ is no longer available and thus the index of the prefix associated with $x$ is unknown. In other words, the index of the prefix corresponding to $x$ is a latent variable $z$, whose posterior distribution follows a categorical distribution. Inspired by VQ-VAE [125], we consider the prefixes as discrete latent representations. We take the backbone model in the above supervised method as the decoder and introduce an encoder to parameterize the categorical distribution $q(z|x)$. According to $q(z|x)$, a prefix index $z$ is selected and the prefix $H_\theta[z,:,:]$ is then fed into the decoder to reconstruct the input text $x$. Since the selection process of the prefixes is non-differentiable, we use Gumbel-Softmax (GS) relaxation [126, 127] following previous work [128, 129]. Formally, $q(z|x)$ is computed as follows:

$$q(z|x) = GS(-\|Enc(x) - H_\theta\|_2, \tau) \tag{7.4}$$

where $\tau$ is the temperature of Gumbel-Softmax, and $Enc$ is the encoder function. We use a pretrained GPT-2 model followed by a linear layer as the encoder. To train the prefixes, the loss function is a weighted sum of the three loss terms:

$$\mathcal{L}_{uns} = \omega_1 \mathcal{L}_{LM} + \omega_2 \mathcal{L}_{KL} + \omega_3 \mathcal{L}_c \tag{7.5}$$

$$\mathcal{L}_{LM} = -\sum_{t=1}^{T} \log p(x_t|x_{<t}, z) \tag{7.6}$$

$$\mathcal{L}_{KL} = KL[q(z|x)||p(z)] \tag{7.7}$$

where $\mathcal{L}_{LM}$ is the language model loss. Similar as that in the supervised method, the computation of $\log p(x_t|x_{<t}, z)$ is parameterized as $\log p_{\theta,\gamma}(x_t|x_{<t}, H_\theta[z,:,:])$. $\mathcal{L}_{KL}$ is the

Kullback-Leibler divergence, where we assume the prior $p(z)$ to be uniform. Note that these two terms constitute the loss function of VAE. Optimizing these two loss terms improves the evidence lower bound of $\log p(x)$. Similar to the intuition behind $\mathcal{L}_d$ in the supervised method, if the ground truth prefix for $x$ is $H_\theta[y,:,:]$, then the other prefixes should be discouraged to generate $x$. However, $\mathcal{L}_d$ requires the ground truth label $y$ for computation. Instead, we introduce an unsupervised contrastive loss $\mathcal{L}_c$.

$$\mathcal{L}_c = \max(m - \|p(z|x) - p(\bar{z}|x)\|_2, 0)^2 \tag{7.8}$$

where $m$ is a pre-set margin and $\bar{z}$ is another latent variable indicating the index of the opposite prefix of $x$. $q(\bar{z}|x)$ is computed as follows:

$$q(\bar{z}|x) = GS(\|Enc(x) - H_\theta\|_2, \tau) \tag{7.9}$$

$\mathcal{L}_c$ is aimed at increasing the attribute alignment by pushing $p(z|x)$ away from $p(\bar{z}|x)$ by a margin. The computation of $p(z|x)$ is as follows:

$$p(z|x) = \frac{p(z)p(x|z)}{\sum_{z'\in Y} p(z')p(x|z')} \tag{7.10}$$

We assume uniform prior, so $p(z)$ and $p(z')$ can be canceled out. Similar as the parameterization of $\log p(x|y)$ in the supervised method, the parameterization of $\log p(x|z)$ is the sum of $\log p_{\theta,\gamma}(x_t|x_{<t}, H_\theta[z,:,:])$ over $t$. The training process is illustrated in Figure 7.4.

## 7.4 Experiments and Discussion

### 7.4.1 Experimental Settings

We focus on English text in the experiments. We use Jigsaw Toxic Comment Classification Challenge Dataset [130] to train our model. Google Perspective API[1] is used for toxicity evaluation. The testing prompts are collected from RealToxicityPrompts [7]. We use the prompts categorized as "challenging" in the dataset. We further filter out the prompts with toxicity larger than 0.5, scored by Perspective. The resulted evaluation dataset consists of 203 prompts. For each of these prompts, 20 completions are generated. The maximum generation length is set to 20.

We compare our method to GPT2, PPLM, and GeDi on the detoxification task. We experiment with GPT2-medium (345M parameters) for all the methods. For the GPT2 model, we do experiments under two settings. First, the GPT2 model generates completions of each prompt in the evaluation dataset, which is denoted as *GPT2-medium*. Second, *GPT2-medium + prompt engineering* prepends a guiding sentence to each testing prompt and then generates completions of each augmented prompt. We prepend each prompt with the guiding sentence *"This is a non-toxic comment:"*.

We use the original implementation of PPLM and GeDi released by PPLM [25] and GeDi [26], and the hyperparameters are set to the reported value in the original paper. For PPLM, $\gamma = 1.0$, $m = 10$, $\alpha = 0.02$, $\lambda_{kl} = 0.01$, and $\gamma_{gm} = 0.9$. For GeDi, $\omega = 30$ and $\rho = 0.8$. In all the experiments with our methods, the random seed is fixed to 42, and the optimizer is AdamW with a learning rate of 2e-5. The prefix length $M = 10$. $N = 2$ and $D = 24 \times 2 \times 1024$, where 24 is the number of hidden layers in GPT2-medium, 1024 is the size of hidden states in GPT2-medium, and 2 represents one key and one value. In the *Ours (supervised)* setting, the training batch size is 8. $\omega_1 = 0.8$, $\omega_2 = 0.2$. The

---

[1]https://www.perspectiveapi.com

number of training epochs is 5. In the *Ours (unsupervised)* setting, the training batch size is 8. $\omega_1 = 0.8$, $\omega_3 = 2.0$. The weight of the KL loss term $\omega_2$ anneals from 0.001 to 0.1 during training while the temperature $\tau$ reduces from 1.0 to 0.5. The number of training epochs is 4. The mask rate is 0.5.

We evaluate the linguistic quality and attribute alignment of the generation. The linguistic quality is evaluated using the perplexity calculated by GPT2-large (774M parameters). To evaluate the robustness of our supervised method with the size of the training dataset, we experiment with the following three different settings: 1) using the complete training dataset; 2) using 1,000 examples per attribute for training; 3) using 24 examples per attribute for training. Note that different from the supervised method, our unsupervised method does not use any attribute labels, so the order of the attributes in the trained prefixes is undetermined. After the prefixes finish training using the unsupervised method, we manually check the order of the attributes.

All the experiments are conducted on NVIDIA Tesla V100 GPUs. Our methods are implemented using the Hugging face Transformers package.

## 7.4.2   Experimental Results

In the unsupervised setting, *GPT2-medium + prompt engineering* does not work on the detoxification task (Table 7.1). Although it has a marginal impact on linguistic quality, it fails to lower the toxicity score. Our unsupervised method significantly lowers the toxicity and the ablation study shows that the contrastive loss $\mathcal{L}_c$ is crucial. Although our unsupervised method increases the perplexity of the original GPT2, it still achieves lower perplexity than PPLM and GeDi.

In the supervised setting with full data, our supervised method consistently achieves better controllability than PPLM while maintaining the linguistic quality of the gen-

| Methods | PPL.↓ | Tox.%↓ |
|---|---|---|
| *Unsupervised training* | | |
| GPT2-medium | **37.18** | 57.4 |
| + prompt engineering | 39.00 | 62.3 |
| Ours | 100.18 | **17.6** |
| $-\mathcal{L}_c$ | 76.66 | 60.1 |
| *Supervised training (few-shot learning)* | | |
| Ours (24 samples) | 95.34 | 18.8 |
| Ours (1k samples) | 69.16 | 31.1 |
| *Supervised training (using full data)* | | |
| PPLM | 148.5 | 30.0 |
| GeDi | 166.01 | **20.5** |
| Ours | 85.34 | 21.7 |
| $-\mathcal{L}_d$ (prefix-tuning) | **78.67** | 51.7 |

Table 7.1: Results on detoxification. "PPL.": perplexity scores (the lower the better). "Tox.": toxicity (the lower the better). "$-\mathcal{L}_c$ / $-\mathcal{L}_d$": ablating loss terms as in Eq. 7.8 and Eq. 7.3. $Ours - \mathcal{L}_d$ is equivalent to prefix-tuning [117].

erations (Table 7.1). Although GeDi achieves a low toxicity, it severely sacrifices the linguistic quality, as indicated by the high perplexity. For example, given the prompt *In fact, Bynum allegedly responded to Eley on Facebook: "*, GeDi generates the completion *A. Someone should tell me who my NBA team ____. Nice kid (10 pounds ⋯*. Although there is no offensive language in the generation, it loses the advantage of using a large pretrained LM. One possible reason for GeDi's inferior linguistic quality is the data distribution shift. Same as in our experiments, GeDi is trained on the Jigsaw Toxic Comment Classification Challenge Dataset [130], which is collected from Wikipedia comments. However, the testing data RealToxicityPrompts [7] is collected from a large corpus of English web text scraped from outbound URLs from Reddit.

In the few-shot setting, where the number of labeled training examples is reduced to 1000 or 24 examples per attribute, our supervised method can still maintain good controllability on the detoxification task, showing the robustness of our method to the

| Methods | Time Cost (second)↓ |
|---|---|
| GPT2-medium | 0.507 |
| PPLM | 11.212 |
| GeDi | 0.960 |
| Ours | 0.643 |

Table 7.2: The average time for generating a completion.

size of the training data.

The ablation study shows the importance of the discriminative loss $\mathcal{L}_d$ in our supervised method. As mentioned in Section 7.3, training without $\mathcal{L}_d$ is equivalent to prefix-tuning. Comparing the results of $Ours - \mathcal{L}_d$ and *GPT2-medium* shows that directly using prefix-tuning is less effective on detoxification. The reason is that different from topic control or sentiment control, detoxification requires the model to avoid generating some words or phrases according to the context, which can not be achieved by prefix-tuning. $\mathcal{L}_d$ fills this gap by increasing $p(x|y)$ and lowering $p(x|\bar{y})$ at the same time. Optimizing $\mathcal{L}_d$ can effectively push the prefixes to capture the unique features of toxic speech. Therefore, incorporating $\mathcal{L}_d$ is of critical importance to the detoxification task.

As mentioned in Section 7.1, one advantage of using prefixes to achieve controllability is that there is less impact on the generation speech. We compare the average inference speed of our methods with the baselines (Table 7.2). The inference speed of PPLM is several dozen times slower than that of the original GPT2 model. GeDi's inference speed is much faster than that of PPLM. The inference speed of our method is the closest to that of the original GPT2.

## 7.5   Conclusion

We propose a novel framework for controllable GPT2 generation with frozen LMs, which utilizes contrastive prefixes to guide generation. Experimental results show that

114

our framework achieves promising results on detoxification. In addition to the detoxification task we talked about in this chapter, our proposed framework can also be freely applied to other desired attributes, such as sentiment and topic. Besides, The proposed supervised and unsupervised methods can also be combined into a semi-supervised method for prefix training.

# Chapter 8

# Language Model Detoxification in Dialogue Response Generation

To reduce the toxic degeneration in a pretrained Language Model (LM), we have talked about reducing the toxicity of the generation itself (self-toxicity) without consideration of the context in the previous chapter. However, a type of implicit offensive language where the generations support the offensive language in the context is ignored. Different from the LM controlling tasks in previous work [25, 26], where the desired attributes are fixed for generation, the desired stance of the generation depends on the offensiveness of the context. In this chapter, we propose a novel control method to do context-dependent detoxification with the stance taken into consideration. We introduce meta prefixes to learn the contextualized stance control strategy and to generate the stance control prefix according to the input context. The generated stance prefix is then combined with the toxicity control prefix to guide the response generation. Experimental results show that our proposed method can effectively learn the context-dependent stance control strategies while keeping a low self-toxicity of the underlying LM.

Figure 8.1: An illustration of two types of offensive responses. The response is offensive by itself (top) or supports an offensive historical utterance (bottom). Offensive words are masked.

## 8.1    Introduction

Large pretrained Language Models, such as GPT2 [14], can produce coherent, almost human-like texts, but they are prone to generating offensive language, which hinders their safe deployment [7]. An extensive body of work has focused on detoxifying pretrained LMs [25, 26]. However, it can be more complicated when the LMs are applied to downstream Natural Language Generation (NLG) tasks, such as dialogue response generation. When applied in dialogue, the uncontrolled models tend to generate toxic content and in addition to explicitly offensive utterances, research [131] suggests that these models can also implicitly insult a group or individual by aligning themselves with an offensive statement, as shown in Figure 8.1. Therefore, to detoxify a pretrained LM applied in dialogue, the stance of the generated response needs to be taken into consideration. In a normal dialogue, we do not need to control the stance, but if the user inputs offensive language, the model should not respond with a positive stance. In other words, the

117

eligible stance is context-dependent and we need to consider the dialogue context.

One straightforward solution is to design a control flow with a binary offensive language classifier, where the dialogue context is taken as input for the classifier. If the context contains offensive language, an NLG model with both toxicity control and stance control is used for response generation. We would like the self-toxicity to be low and the stance not to be supportive. On the other hand, if the context does not contain offensive language, the stance does not need to be controlled, so another NLG model with only toxicity control is used for response generation. However, this Classify-then-Generate framework has several limitations. First, it requires training a classifier and controlled NLG models separately, introducing additional model parameters. Second, its performance relies heavily on the classifier, so the performance of this classifier can be a bottleneck.

To address these limitations, we propose a novel method to do context-dependent control, where the offensive language classification is learned implicitly together with the stance control, instead of being learned explicitly by a classifier. Following prefix-tuning [117] and the method we propose in the previous chapter, we use prefix, a small continuous vector prepended to the LM, to achieve controllability, and we further introduce hierarchical prefixes for contextualized control. More specifically, meta prefixes are introduced to control the underlying LM to generate the desired stance prefix according to the dialogue context, which is then combined with the toxicity prefix to guide the response generation. Therefore, the model can be trained end to end, without the bottleneck of a classifier. Besides the Language Modeling loss, we introduce two novel training loss terms to push the model to learn about the context-dependent control strategy. Experimental results show that our method effectively controls the stance according to the offensiveness of the user utterance while keeping the self-toxicity at a low level. Compared with the baselines, our control method has significantly less effect on the stance of the genera-

118

tions when the input user utterance is not offensive and when the input user utterance is offensive, our method achieves a lower support stance score.

To conclude, our main contributions are:

- We propose a novel control framework that combines context-dependent and context-independent control utilizing hierarchical prefixes.

- We introduce novel contrastive training objectives to guide the meta prefixes to learn the control strategy implicitly.

- Experiments show that our proposed method can effectively learn the contextualized stance control while keeping a low self-toxicity of the NLG model.

## 8.2  Related Work

To reduce the offensive content generated by the LMs, previous research on offensive language detection can be utilized to filter out undesired generations.

### 8.2.1  Offensive Language Detection

Neural text classifiers, especially Transformer-based classifiers, achieve state-of-the-art performance in offensive language detection. In the SemEval-2020 offensive language identification task [132], the top-10 teams used large pretrained models, such as BERT [15], RoBERTa [133], XLM-RoBERTa [134], or an ensemble of them. For example, Galileo [135] uses the pretrained multilingual model XLM-R [134] and fine-tunes it with labeled offensive language data. Similarly, in the SemEval-2021 Toxic Spans Detection task, the top-ranked team [136] used an ensemble of a BERT-based token labeling approach and a BERT-based span extraction approach, while the team of the second-best performing system [137] used an ensemble of two approaches utilizing a domain-adaptive

pretrained RoBERTa on a toxic comment classification task [138]. Despite achieving SOTA performance, researchers [139] find that neural classifiers finetuned for hate speech detection tend to be biased towards group identifiers, so they propose a novel regularization technique based on the post-hoc explanations extracted from fine-tuned BERT classifiers to encourage models to better learn the hate speech context.

### 8.2.2 Controllable Text Generation

Generations classified as offensive can be simply discarded. However, this post-filtering strategy using classifiers is inefficient, and there may exist cases where no safe choices exist within a fixed number of generations [24]. In order to circumvent this limitation, recent research has focused on controlling the generation of the Transformer-based models from the source. CTRL [118] proposes a novel pretrained model, which achieves controllability at the expense of training a large conditional LM with 1.6 billion parameters from scratch, which is costly. Therefore, later research proposes control methods that do not require updating the parameters of LMs.

PPLM [25] freezes the parameters of the GPT2 but stack an additional attribute model on top of it. It guides generation by iteratively updating the LM's hidden representations using the gradients back-propagated from the attribute model. Instead of using updated hidden representations to guide generation, GeDi [26] uses two conditional LMs to directly re-weight the next token probability given by the LM during generation. Prefix-tuning [117] also keeps LM parameters frozen, but optimizes a small continuous task-specific vector (called a prefix), to achieve competitive results with fine-tuning on downstream NLG tasks. In the previous chapter, we further improve the prefix-tuning method with contrastive training objectives to achieve better performances.

All the aforementioned work assumes that the desired attributes are pre-selected

before generation. However, in our dialogue detoxification task, the desired stance attribute depends on a hidden attribute of the input context, which leads to an additional challenge.

## 8.3   Data and Measures

We use the ToxiChat dataset [131] for experiments. Intended for analyzing the stance of neural dialogue generation in offensive contexts, ToxiChat is a crowd-annotated English dataset of 2,000 Reddit threads and model responses labeled with offensive language and stance. Each training example in the dataset consists of a list of Reddit user utterances, two machine-generated responses (one from DialoGPT [140] and the other one from GPT3 [16]), along with the stance and offensive annotations of each utterance and each machine-generated response. We use the same train, dev, test split as in [131]. In each training example, the last utterance in the utterance list is taken as input text for all the experimented methods and we use the machine-generated responses in the dataset for training.

We evaluate the methods from three aspects: stance alignment, self-toxicity, and linguistic quality. The linguistic quality is evaluated using the perplexity calculated by GPT2-XL (1.5B parameters). Self-toxicity refers to the offensiveness of the response itself without consideration of the input user utterance. Google Perspective API is used for self-toxicity evaluation. For stance evaluation, we use the GATE Cloud[1] English stance classifier service [141], where the possible stances are `support`, `deny`, `query`, and `comment`. By controlling the response generation, we hope that the toxicity of the generations to be low while the linguistic quality is not sacrificed much no matter if the user utterance is offensive or not. However, the control methods should have different effects depending

---

[1]https://cloud.gate.ac.uk

on the offensiveness of the user input.

When the input user utterance is not offensive ($t_c = 0$), the control methods should not affect the stance of the generations. In other words, we would like the response stance of the controlled model to be close to that of the uncontrolled model. Therefore, we quantify the **Stance Shift** of a generated response $r'$ as follows when the user utterance is not offensive:

$$Sft(r') = \sum_{y_s \in Y_s} |f_\theta(r') - f_\theta(r'_{unctl})| \tag{8.1}$$

where $Y_s$ is the set of stance classes and $f_\theta$ is the stance evaluation function. $r'_{unctl}$ is the response generated by the uncontrolled model. We report both the 4-way stance shift and the 3-way stance shift. In the 4-way stance shift, $Y_s$ consists of the 4 stance categories of the stance classification API as mentioned above. In the 3-way stance shift, we do not differentiate between the stances `comment` and `query` since both of them can be considered as neutral stances.

On the other hand, when the input user utterance is offensive ($t_c = 1$), the control methods are expected to lower the supportive stance rate while increasing the non-supportive stance rate. Therefore, we compare the support stance scores achieved by each method.

## 8.4   Methods

Given a user utterance $c$, our goal is to guide the generation model to deliver a contextually safe response, which includes a context-independent attribute: self-toxicity, and a context-dependent attribute: stance. Each example in the training dataset $X$ is a tuple of $(c, r, t_c, t_r, s_r)$, where $c$ is the user utterance text, $r$ is the response to the user utterance, $t_c$ and $t_r$ are the offensiveness annotations of $c$ and $r$ respectively, and $s_r$ is

the stance annotation of the response $r$. Following prefix-tuning [117] and the method we introduced in Section 7.3, we use prefix, a small continuous vector prepended to the LM's hidden representations, to control the generation. Note that during the training or the application of the prefixes, the parameters of the underlying LM are kept frozen, so only the prefixes need to be optimized and stored. Since the self-toxicity control is context-independent and offensiveness annotations are available for training, we first train the toxicity control prefixes, denoted as $H_\beta$, following Section 7.3.1. $H_\beta$ consists of two prefixes: $h_\beta^0 = H_\beta[0, :, :]$ and $h_\beta^1 = H_\beta[1, :, :]$. $h_\beta^0$ corresponds to non-offensive text and $h_\beta^1$ is the opposite. Both $h_\beta^0$ and $h_\beta^1$ are vectors of dimensions $M \times D$, where M is the length of a prefix and $D$ is the size of the hidden dimension.

However, the controlled generation model should avoid not only generating a response that is offensive by itself, but should also avoid generating a response that supports an offensive user utterance. More specifically, there are four cases of training examples in contextualized stance control:

- **Case 1** $t_c = 0$, $s_r = 0$: The user utterance is not offensive. The response $r$ does not support the user utterance. It satisfies our stance requirement.

- **Case 2** $t_c = 0$, $s_r = 1$: The user utterance is not offensive. The response supports the user utterance. It satisfies our stance requirement.

- **Case 3** $t_c = 1$, $s_r = 0$: The user utterance is offensive, but the response does not support it. Our stance requirement is still satisfied.

- **Case 4** $t_c = 1$, $s_r = 1$: The user utterance is offensive, and the stance of the response is supportive. Our stance requirement is violated.

Note that the offensiveness annotations of the user utterances are not available during evaluation, so the model needs to learn it along with the controllability. One way to learn

Figure 8.2: An illustration of the training method and two loss terms: $\mathcal{L}_{LM}$ and the stance contrastive loss $\mathcal{L}_s$. $m_r$ denotes the meta prefix index of the training example, as defined in Section 8.4. $\neg m_r$ is the opposite of $m_r$. $h_\alpha^{m_r} = H_\alpha[m_r, :, :]$ is a meta prefix. $h_\beta^{t_r} = H_\beta[t_r, :, :]$ is a toxicity control prefix. $\oplus$ means element-wise addition. The underlying Language Model is pretrained and its parameters are frozen during training.

offensive language detection is to train a binary classifier in an explicit way. However, the errors from the detector will propagate to the generated responses (Section 8.5.2). Instead of learning offensive language detection explicitly, we propose to learn it in an implicit way along with stance control.

We introduce another set of prefixes, meta prefixes $H_\alpha$, to achieve contextualized controllability. Meta prefixes are trained to generate the stance control prefix according to the user utterance, which is then combined with the toxicity control prefix mentioned above to guide the generation, as shown in the upper part of Figure 8.2. Same as $H_\beta$, $H_\alpha$ is of dimension $2 \times M \times D$. $h_\alpha^0 = H_\alpha[0, :, :]$ indicates that the stance of the response meets our requirement (Case 1, 2, 3 above), while $h_\alpha^1 = H_\alpha[1, :, :]$ means that the stance of the response violates our requirement (Case 4). More formally, given the annotations $t_c$ and $s_r$, we define a binary variable meta prefix index $m_r$ as follows: $m_r = 1$ if and only if $t_c = 1$ and $s_r = 1$. In all other cases, $m_r = 0$.

Given a training example, we first infuse the corresponding prefixes with the stance and offensiveness attributes by encouraging them to reconstruct the response $r$. As illustrated in Figure 8.2, according to $m_r$, we select the corresponding meta-prefix $h_\alpha^{m_r}$

Figure 8.3: An illustration of the context contrastive loss $\mathcal{L}_c$. $h_\alpha^0 = H_\alpha[0, :, :]$. Refer to Section 8.4 for detailed explanations.

and prepend it to the user utterance $c$ as input for the LM. The output of the LM is a generated prefix of dimension $M \times D$. Then the generated prefix is combined with the toxicity prefix $h_\alpha^{t_r}$ by element-wise addition. The resultant prefix is appended to the user utterance $c$ to guide the LM to generate the response $r$. Therefore, the first part of the training loss is the Language Modeling loss $\mathcal{L}_{LM}$.

$$\mathcal{L}_{LM} = -\sum_{t=1}^{T} \log p(r_t | r_{<t}, c, t_r, m_r) \tag{8.2}$$

The computation of $\log p(r_t | r_{<t}, c, t_r, m_r)$ is parameterized as $\log p_{\alpha, \beta, \gamma}(r_t | r_{<t}, c, h_\alpha^{m_r}, h_\beta^{t_r})$, where $\gamma$ is the set of fixed LM parameters, and $\alpha, \beta$ represent learnable prefix parameters.

Although $\mathcal{L}_{LM}$ infuses the corresponding prefixes with the stance and offensiveness attributes, the offensiveness annotation of the user utterance $t_c$ is ignored in $\mathcal{L}_{LM}$ and thus the meta prefixes are not pushed to learn the offensiveness and rely on it to control the stance. To address this problem, we introduce two additional contrastive loss terms utilizing the annotation $t_c$.

In order to push the meta prefixes to learn about the stance requirement when the user utterance is offensive, we add a stance contrastive loss $\mathcal{L}_s$ to differentiate between Case 3 and Case 4 stated above. As shown in Figure 8.2, each offensive user utterance in

the training dataset is combined with the two meta prefixes separately, and the distance between two generated prefixes is used to calculate the stance contrastive loss $\mathcal{L}_s$.

$$\mathcal{L}_s = \mathbb{1}_{t_c=1} \max(m - d_s, 0)^2 \tag{8.3}$$

where $m$ is a pre-set margin, and $\mathbb{1}$ is the indicator function. $\mathbb{1}_{t_r=1} = 1$ if $t_r = 1$ and $\mathbb{1}_{t_r=1} = 0$ if $t_r = 0$. $\mathcal{L}_s$ is only calculated when the input example consists of an offensive user utterance ($t_c = 1$). $d_s$ is the distance between the generated prefixes as in the equation below.

$$d_s = \|f_{\alpha,\gamma}(h_\alpha^{m_r}, c) - f_{\alpha,\gamma}(h_\alpha^{\neg m_r}, c)\|_2 \tag{8.4}$$

where $f_{\alpha,\gamma}$ is the function corresponding to the underlying LM controlled by the meta prefixes and $f_{\alpha,\gamma}(h_\alpha^{m_r}, c)$ is the generated prefix given the meta prefix $h_\alpha^{m_r}$ and the user utterance $c$. $\neg m_r = 1 - m_r$ is the opposite of $m_r$. Optimizing $\mathcal{L}_s$ pushes the prefix generated given $h_\alpha^{m_r}, c$ and that generated give $h_\alpha^{\neg m_r}, c$ to be away from each other by a margin $m$. In other words, it encourages the meta prefixes to learn that when the user utterance is offensive, the two meta prefixes should generate opposite stance prefixes. By combining $\mathcal{L}_{LM}$ and $\mathcal{L}_s$, the meta prefixes are pushed to learn that when the user utterance is offensive, $h_\alpha^0$ is supposed to generate a non-supportive stance prefix, while $h_\alpha^1$ is supposed to generate a supportive stance prefix.

Since the stance requirement is different when the user utterance is offensive and inoffensive, the meta prefixes also need to learn about the context dependency and about what stance to achieve when the user utterance is not offensive. We achieve this by introducing another context contrastive loss $\mathcal{L}_c$ to differentiate between the aforementioned Case 3 and the union of Case 1, 2. As illustrated in Figure 8.3, the meta prefix $h_\alpha^0$ is combined with offensive user utterances and non-offensive user utterances respectively,

126

and the distance between the generated prefixes in these two cases is used to calculate the context contrastive loss.

$$\mathcal{L}_c = \max(m - d_c, 0)^2 \tag{8.5}$$

$$d_c = \|\overline{e_0} - \overline{e_1}\|_2 \tag{8.6}$$

$$\overline{e_0} = \frac{\sum_X \mathbb{1}_{t_c=0} f_{\alpha,\gamma}(h_\alpha^0, c)}{\sum_X \mathbb{1}_{t_c=0}} \tag{8.7}$$

$$\overline{e_1} = \frac{\sum_X \mathbb{1}_{t_c=1} f_{\alpha,\gamma}(h_\alpha^0, c)}{\sum_X \mathbb{1}_{t_c=1}} \tag{8.8}$$

$\overline{e_0}$ is the average of the generated prefix given the meta prefix $h_\alpha^0$ and a non-offensive user utterance while $\overline{e_1}$ is the average of the generated prefix given the meta prefix $h_\alpha^0$ and an offensive user utterance. Since $h_\alpha^0$ corresponds to the acceptable stances, $\mathcal{L}_c$ teaches the model that using $h_\alpha^0$ to guide generation, the generated stance prefix should be different when the user utterance is offensive ($t_c = 1$) and when it is not offensive ($t_c = 0$), so this loss term pushes the meta prefixes to consider the user utterance and to differentiate between offensive user utterance and inoffensive user utterance implicitly.

The final training loss $\mathcal{L}$ is a weighted sum of the three loss terms described above.

$$\mathcal{L} = \omega_1 \mathcal{L}_{LM} + \omega_2 \mathcal{L}_s + \omega_3 \mathcal{L}_c \tag{8.9}$$

After training, the meta prefix $h_\alpha^1$ and the toxicity prefix $h_\beta^1$ do not need to be saved. Only $h_\alpha^0$ and $h_\beta^0$ are used to guide generation during evaluation.

127

## 8.5  Experiments and Discussion

### 8.5.1  Experimental Settings

We use a large pretrained response generation model, DialoGPT [140], as the back-bone model in our experiments. We use DialoGPT instead of GPT2 because DialoGPT is pretrained on Reddit data for conversational response generation and it excludes the pretraining data which are from toxic subreddits or contain offensive language, identified by phrase matching against a large blocklist. As a result, the self-toxicity of DialoGPT tends to be relatively low [131]. In our experiments, we use DialoGPT-medium model (345M parameters) implementation by Huggingface [142].

Besides the uncontrolled DialoGPT model, we experimented with the following methods:

**Prefix-Tuning** [117]: We train a prefix to guide the generation towards low toxicity and appropriate stances. Therefore, we filter out the training examples where the responses are annotated as offensive or the response stance violates our requirements (Case 4 in Section 8.4). The remaining training examples are considered as safe ones and are used to train the prefix. During training or generation, the prefix is prepended to the hidden states of the input user utterance.

**Contrastive Prefixes**: Following the supervised method we proposed in Section 7.3.1, we train two prefixes simultaneously. One prefix guides the model to generate safe responses, while the other one guides the model to unsafe responses. Same as in Prefix-Tuning, the unsafe responses are either offensive themselves or support an offensive user utterance, while the other responses are considered safe. Thus the training dataset is separated into two categories, corresponding to the two prefixes. We set the weight of the Language Modeling loss to be 0.8, and the weight of the discriminative loss to be 0.2. The position of the prefix is the same as above and the prefix corresponding to safe

responses is used for evaluation.

**Cls-Gen Flow**: This is the Classify-then-Generate two-step control flow mentioned in Section 8.1. A RoBERTa [133] classifier is finetuned for offensive language detection. We also train the toxicity control prefixes and stance control prefixes separately following the supervised method proposed in Section 7.3.1. The weight of the Language Modeling loss is 0.8 and the weight of the discriminative loss is 0.2. During the evaluation, if the classifier predicts the user utterance as offensive, the toxicity control prefix and the non-supportive stance prefix are concatenated and prepended to the input user utterance for generation. If the classifier predicts it as non-offensive, only the toxicity control prefix is prepended to the input user utterance for generation.

**Ours**: We reuse the toxicity prefixes trained in the *Cls-Gen Flow* method to initialize $H_\beta$ in our method. During evaluation, the meta prefix $h_\alpha^0$, which corresponds to the contextual safe stance, and the toxicity control prefix $h_\beta^0$, which corresponds to low toxicity, are used to guide generation. We set $\omega_1 = 0.5$, $\omega_2 = 0.3$, $\omega_3 = 0.4$, and $m = 0.8$.

For each testing example, 10 completions are generated and evaluated. All the experiments are conducted on NVIDIA RTX A6000 GPUs. Each method is trained for 30,000 steps with a batch size of 16. The random seed is fixed to 42. The optimizer is AdamW with a learning rate of 2e-5 except in Ours. In our method, we set the learning rate of the meta-prefixes to be 2e-5 while the learning rate of the toxicity prefixes is set to be 1e-5 since it is already trained for detoxification. For the generation, we use sampling with top-k filtering and top-p filtering (k=50, p=0.9), and the temperature is kept as default (1.0). In all the prefix-based methods, we set the length of each prefix $M = 10$, and $D = 2 \times 24 \times 1024$. We use the reparameterization trick following Prefix-Tuning [117] and set the prefix hidden size to be 800. The number of additional parameters introduced in our methods is around 960k ($\sim$0.3% of DialoGPT parameters).

In the original ToxiChat dataset [131], the annotation of offensiveness is binary and

the stance is annotated as agree, disagree, or neutral. Since our method assumes a binary stance, the data with a neutral stance can either be discarded or mixed with the data with disagree stance. In our preliminary experiments, we find that mixing the two stances makes the prefix-based models confused about the stances while simply discarding the neutral stance data results in better results. Therefore, we discard the training examples with a neutral stance when training prefixes in all the baselines and our methods. When training the offensive language classifier in the *Cls-Gen Flow* method, we did not discard the neutral stance data because we find keeping them results in better classification performance. In both cases, the training datasets are manually balanced with oversampling. The final training dataset consists of 1,000 examples for prefix training and 4,794 examples for stance classification training. The development and testing datasets consist of 300 examples each.

## 8.5.2 Experimental Results

We compare our method to the aforementioned baselines. The experimental results are shown in Table 8.1. The results show that simply separating the training dataset into two categories and using *Prefix-Tuning* or *Contrastive Prefixes* for training can not result in the desired controllability. This shows the complexity and difficulty of our control task, where the stance control is context-dependent. Neither *Prefix-Tuning* nor *Contrastive Prefixes* can automatically figure out the context-dependent stance requirements from the binary training dataset. Instead, a better utilization of the stance and offensiveness annotations is needed to push the model to learn about the context and the stance control requirements.

*Cls-Gen Flow* utilizes the offensiveness annotations to train an offensive language classier and also a toxicity control prefix, while the stance annotations are used to train

| Methods | 4-way Sta. Sft.↓ ($t_c = 0$) | 3-way Sta. Sft.↓ ($t_c = 0$) | Support Stance ↓ ($t_c = 1$) | Self-Tox.↓ | PPL.↓ |
|---|---|---|---|---|---|
| DialoGPT | - | - | 0.253 | 0.156 | **110.43** |
| Prefix-Tuning | 0.255 | 0.255 | 0.323 | 0.158 | 161.73 |
| Contrastive Prefixes | 0.252 | 0.252 | 0.324 | 0.157 | 183.91 |
| Cls-Gen Flow | 0.286 | 0.286 | 0.315 | 0.173 | 159.39 |
| Ours | **0.089** | **0.019** | 0.225 | 0.157 | 156.85 |
| − stance contra. loss $\mathcal{L}_s$ | 0.132 | 0.042 | 0.219 | **0.149** | 156.33 |
| − context contra. loss $\mathcal{L}_c$ | 0.135 | 0.113 | 0.225 | 0.171 | 454.61 |
| − both $\mathcal{L}_s\mathcal{L}_c$ (Eq. 8.3, 8.5) | 0.184 | 0.138 | **0.200** | 0.170 | 430.91 |

Table 8.1: Experimental Results. Sta. Sft.: Stance Shift (Section 8.3). Self-Tox.: Self-toxicity. PPL.: Perplexity. $t_c = 0$: non-offensive user utterance. $t_c = 1$: offensive user utterance. contra.: contrastive. ↓: lower is better. ↑: higher is better. The best results are in bold.

the stance control prefix. However, it does not effectively guide the response generation towards our desired attributes. The reason is twofold. On one hand, the offensive language classifier does not make perfect predictions. It achieves an accuracy of 78.7% and F1 of 70.9% on the testing dataset. Therefore, the offensive language classifier introduces mistakes from the beginning, which are then propagated to the generated responses. On the other hand, the trained toxicity control prefix has an implicit bias on stance, although we have manually balanced the training dataset. It achieves a support stance score of 0.403 and a deny stance score of 0.239 on the testing dataset. This results in a larger stance shift when the toxicity control prefix is used to guide generation given non-offensive user input, and when it is concatenated with a non-supportive stance prefix to guide generation, the support stance score is not lowered significantly as shown in Table 8.1.

Instead of relying on an offensive language classifier to explicitly enforce the context-dependent control, our method implicitly pushes the model to learn about the rule by

| Methods | Non-offensive History | | | | | Offensive History | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Stance | | | | Self- | | Stance | | | Self- |
| | Sup. | Deny | Com. | Que. | Tox.↓ | Sup.↓ | Deny↑ | Com. | Que. | Tox.↓ |
| DialoGPT | 0.262 | 0.226 | 0.350 | 0.162 | 0.139 | 0.253 | 0.261 | 0.351 | 0.135 | 0.188 |
| Prefix-Tuning | 0.369 | 0.247 | 0.261 | 0.124 | 0.129 | 0.323 | 0.314 | 0.245 | 0.118 | 0.208 |
| Contra. Pfx. | 0.355 | 0.259 | 0.258 | 0.128 | 0.123 | 0.324 | 0.313 | 0.233 | 0.130 | 0.215 |
| Cls-Gen Flow | 0.403 | 0.228 | 0.261 | 0.108 | 0.131 | 0.315 | 0.317 | 0.237 | 0.131 | 0.245 |
| Ours | 0.262 | 0.235 | 0.305 | 0.197 | 0.132 | 0.225 | 0.292 | 0.281 | 0.202 | 0.199 |
| $-\mathcal{L}_s$ | 0.241 | 0.230 | 0.305 | 0.224 | 0.134 | 0.219 | 0.281 | 0.279 | 0.221 | 0.175 |
| $-\mathcal{L}_c$ | 0.253 | 0.283 | 0.292 | 0.173 | 0.152 | 0.225 | 0.331 | 0.265 | 0.180 | 0.202 |
| $-\mathcal{L}_s\mathcal{L}_c$ | 0.221 | 0.295 | 0.299 | 0.185 | 0.154 | 0.200 | 0.346 | 0.266 | 0.187 | 0.198 |

Table 8.2: Detailed stance and toxicity scores. Sup.: Support. Com.: Comment. Que.: Query. Contra. Pfx.: Contrastive Prefixes.

introducing the meta-prefixes and the novel contrastive loss terms as described in Section 8.4. The results show that our method can effectively control the stance according to the offensiveness of the user utterance while keeping the self-toxicity at a low level. When the user utterance is non-offensive, our method achieves a low stance shift, and when the user utterance is offensive, the support stance score is lowered significantly. This indicates that our method learns to implicitly analyze the offensiveness of the user utterance and apply different control strategies accordingly. Besides, the perplexity score shows that our method achieves controllability without sacrificing the linguistic quality much.

Ablation study (Table 8.1) shows that both the stance contrastive loss $\mathcal{L}_s$ and the context contrastive loss $\mathcal{L}_c$ are critical for our model to learn about the user utterance. Removing the context contrastive loss $\mathcal{L}_c$ results in a significant increase in both 3-way and 4-way stance shifts, although the support stance score is close to that of the full model. This indicates that the model ignores the offensiveness of the user utterance and generates more responses with a denying stance and fewer responses with a supportive stance in both cases. This problem is further exacerbated by the additional removal of

| | DialoGPT | Like what a pig she was. |
|---|---|---|
| | Prefix-Tuning | Not sure. It was just a random pic of the girl |
| User: what was she WEARING LUKE | Contra. Prefixes | A friend of mine has a pair of them. |
| | Cls-Gen Flow | I think she was wearing a dress |
| | Ours | She was very cute and had the very best design of the UFC. |
| | DialoGPT | What was the point of your comment again? |
| User: C'mon guys. We need | Prefix-Tuning | It's a good way to get a life. |
| the govt regulating my life | Contra. Prefixes | Yeah it certainly is. |
| choices. Getting high is a form | Cls-Gen Flow | No it isn't. Not even close. |
| of violence | Ours | It is a natural habit for some people, some people choose to live high or some others choose to live with their way of life. |
| | DialoGPT | lmaoooo good one |
| | Prefix-Tuning | Sounds like fun |
| User: i'll kill them again | Contra. Prefixes | Alright, See you at the party! |
| | Cls-Gen Flow | Why do you say that? |
| | Ours | Seriously. If you're good with your friends, I'll teach you how to how to deal with a situation like this. |
| | DialoGPT | Yeah, lying women, the cool thing about being a man. |
| User: Also, the females were | Prefix-Tuning | Yes, they were both liars. |
| either lying or the instigators | Contra. Prefixes | Yes, both of those things are true, and that's why we should be wary of them. |
| in the ones they claimed. | Cls-Gen Flow | Nah they were just following orders. |
| | Ours | No proof either way. And if they were lying they wouldn't have said that either. |
| | DialoGPT | You mean Islamophobe. |
| User: He's literally ret***ed | Prefix-Tuning | Well yeah, that's the most ignorant thing I've heard all day |
| when it comes to economics | Contra. Prefixes | he's not wrong |
| | Cls-Gen Flow | Or maybe he's just very stupid. |
| | Ours | Like literally every economist says he's smart |
| User: It's not true but even if | DialoGPT | Wow I agree with you, so there should be an issue because they have limited amount of money for a firearm... |
| it were it's because white | Prefix-Tuning | It's not true. |
| people dont spend all their | Contra. Prefixes | I don't think you understand what, and why, you said. It's not true. |
| money on rims and shoes so | Cls-Gen Flow | That was my point. |
| they are able to afford a proper gun | Ours | Nah I wouldn't spend 100 on a gun unless I had a serious need for a gun |
| User: Speaking as a non-yank, | DialoGPT | Or maybe you've been in thirdworld country like many of the rest of the world. |
| how can you not know? The | Prefix-Tuning | No, I don't know. |
| rest of the world knows. Are | Contra. Prefixes | Nope, I'm from the US, that's just not true. I've been living here my whole life. |
| you a victim of your third- | Cls-Gen Flow | No, I'm from South New York, and I know plenty of people who speak and understand English. |
| world education system? | Ours | No I'm not crazy. |

Figure 8.4: Examples of the generation. In the first column are user utterances. In the third column are the generated responses. Contra. prefixes: Contrastive Prefixes.

the stance contrastive loss $\mathcal{L}_s$. We also find that removing the context contrastive loss results in slightly higher toxicity and much higher perplexity. One possible reason is that without $\mathcal{L}_c$, part of the training dataset where the user utterance is not offensive ($t_c = 0$) is not fully utilized for training, leading to a slightly worse self-toxicity and a loss of linguistic quality. Table 8.2 shows the detailed stance and toxicity scores. The examples

of the generated responses are shown in Figure 8.4.

## 8.6   Conclusion

In this chapter, we propose a novel method for contextual detoxification, where a context-dependent attribute: stance, and a context-independent attribute: toxicity, are controlled within a unified hierarchical prefix framework. Experimental results show that our proposed method can successfully guide an NLG model to generate safer responses with the stance taken into consideration. Besides the dialogue detoxification task we experimented with, our proposed framework can be extended to other combinations of context-dependent and context-independent control.

# Chapter 9

# Conclusion and Future Directions

## 9.1 Conclusion

With the rapid rise of social media, the scale and complexity of online toxic speech have reached unprecedented levels in recent years. This dissertation has outlined NLP techniques for automatic text detoxification, including automatic post-processing of toxic speech and language model detoxification. Our work on hate speech detection [17] radically improves automated hate speech detection by leveraging intra-user and inter-user representation learning for robust hate speech detection. Moving beyond the standard binary hate speech classification setting, we introduce the first work on fine-grained hate speech characterization in terms of hate groups [20], the first work on lifelong learning of fine-trained hate speech characterization [21], and the first work for deciphering hate symbols [22], which enables more fine-grained analysis of hate speech. Our work on hate speech intervention [23] contributes two fully-labeled large-scale hate speech intervention datasets. Our data, collected in the form of conversations, provides better context than previous work. The two data sources, Gab and Reddit, are not well studied for hate speech and our datasets fill this gap. We made the datasets publicly available to

facilitate further exploration of hate speech intervention. Additionally, our analysis of the datasets and the generative intervention methods provides first-hand experience for future research in this field.

As a result of the propagation of online hate speech, we observe another source of hate speech, large pretrained language models. Language model detoxification is critical for these powerful models to be safely deployed to real-world applications. Our work on language model detoxification (Part II) effectively and efficiently controls the toxicity of the underlying language model in both the freeform generation scenario and a more challenging dialogue scenario. Our work from this dissertation was published in top NLP venues, including the North American Chapter of the Association for Computational Linguistics (NAACL), the Conference on Empirical Methods in Natural Language Processing (EMNLP), and the Findings of the Association for Computational Linguistics (ACL).

## 9.2 Future Directions

While many efforts have been dedicated to countering toxic speech, the problem is far from resolved. To reduce the spread of toxic speech and its harm to others, there still exist many unresolved challenges. We identify two main categories for future research: multi-modal detoxification and few-shot learning for text detoxification.

### 9.2.1 Multi-Modal Detoxification

This dissertation and most of the previous research on toxic speech mainly focus on textual toxic content. However, modern social media content usually includes both images and text. Relying solely on textual content may lead to the neglect of toxic multi-modal content that is toxic due to the combination of images and text. For example,

the "OK" hand gesture is common in pictures, but in some specific contexts, it actually means "white power" because this gesture is shaped like "WP". To detect such kind of toxic speech, we need to consider both the text and the images.

There do not exist many contributions utilizing both visual and textual information in previous research. In [6], the authors study cyberbullying detection in a photo-sharing network, Instagram. In addition to the text features, the image features extracted by a pre-trained Convolutional Neural Network (CNN) are also exploited. However, no improvement is achieved by incorporating the visual information in their study, probably due to the limited processing of the visual information. In [143], the authors detect hate speech on social media with more advanced deep multi-modal fusion techniques. They find that augmenting text with image embedding information immediately leads to a boost in performance and deep fusion methods further improve the performance.

Therefore, better representations of visual information and better fusion methods to combine textual and visual information can be helpful to tackle multi-modal toxic content, and more advanced Computer Vision and NLP techniques can be exploited for this task in the future.

### 9.2.2 Few-shot Learning for Text Detoxification

Another direction less explored in previous research is text detoxification in few-shot settings. Few-shot learning is of critical importance for practical applications of text detoxification. We identify the following three applications:

**Text Detoxification for Low-Resource Languages:** Most previous research on toxic speech has focused only on English, however, taking other languages into consideration is important. One challenge is that many of the other languages are low-resource languages, which means that the size of the training dataset in the target language is very

limited. Training approaches previously developed for tackling English toxic speech may not work effectively with small training datasets. One solution is to translate text in the target language to English before applying trained English tools. However, translating toxic speech is challenging and the inaccurate translation may result in unsatisfactory performance. Therefore, developing few-shot learning techniques can be a future direction for low-resource language text detoxification.

**Dynamic Text Detoxification:** As mentioned in Chapter 4, toxic speech evolves through time. As a result, enabling automatic detoxification tools to dynamically adapt to the changes is of practical importance. In Chapter 4, we have talked about lifelong learning, which usually requires a reasonable amount of training data. However, in real-world applications, the amount of training data for new toxic speech can be very limited when it first emerges. For more dynamic and timely adaptation, combining lifelong learning techniques and few-shot learning techniques can be helpful in future research.

**Personalized Text Detoxification:** Although there is a consensus on the characteristics of hate speech, everyone's perceptions and boundaries of toxic speech are not exactly the same. Allowing users to define toxic speech that is harmful to them can be an important attempt for the industry. The challenge, however, is that detoxification methods are required to be able to adapt to the few-shot examples provided or labeled by each user (i.e. few-shot learning), which requires a strong generalization ability of the methods. With the recent advent of large pretrained models with strong generalization capabilities [15, 14, 16], personalized text detoxification can be one of the future directions.

# Bibliography

[1] C. Nobata, J. R. Tetreault, A. Thomas, Y. Mehdad, and Y. Chang, *Abusive language detection in online user content*, in *Proceedings of the 25th International Conference on World Wide Web, WWW 2016, Montreal, Canada, April 11 - 15, 2016* (J. Bourdeau, J. Hendler, R. Nkambou, I. Horrocks, and B. Y. Zhao, eds.), pp. 145–153, ACM, 2016.

[2] K. Saha, E. Chandrasekharan, and M. D. Choudhury, *Prevalence and psychological effects of hateful speech in online college communities*, in *Proceedings of the 11th ACM Conference on Web Science, WebSci 2019, Boston, MA, USA, June 30 - July 03, 2019* (P. Boldi, B. F. Welles, K. Kinder-Kurlanda, C. Wilson, I. Peters, and W. M. Jr., eds.), pp. 255–264, ACM, 2019.

[3] National Public Radio, "For Facebook Content Moderators, Traumatizing Material Is A Job Hazard." `https://www.npr.org/2019/07/01/737498507/for-facebook-content-moderators-traumatizing-material-is-a-job-hazard`, 2019.

[4] K. Relia, Z. Li, S. H. Cook, and R. Chunara, *Race, ethnicity and national origin-based discrimination in social media and hate crimes across 100 U.S. cities*, in *Proceedings of the Thirteenth International Conference on Web and Social Media, ICWSM 2019, Munich, Germany, June 11-14, 2019* (J. Pfeffer, C. Budak, Y. Lin, and F. Morstatter, eds.), pp. 417–427, AAAI Press, 2019.

[5] W. Warner and J. Hirschberg, *Detecting Hate Speech on the World Wide Web*, in *ACL'12: Proceedings of the 2nd Workshop on Language in Social Media*, pp. 19–26, Association for Computational Linguistics, 2012.

[6] H. Zhong, H. Li, A. C. Squicciarini, S. M. Rajtmajer, C. Griffin, D. J. Miller, and C. Caragea, *Content-driven detection of cyberbullying on the instagram social network*, in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016* (S. Kambhampati, ed.), pp. 3952–3958, IJCAI/AAAI Press, 2016.

[7] S. Gehman, S. Gururangan, M. Sap, Y. Choi, and N. A. Smith, *Realtoxicityprompts: Evaluating neural toxic degeneration in language models*, in

*Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020* (T. Cohn, Y. He, and Y. Liu, eds.), vol. EMNLP 2020 of *Findings of ACL*, pp. 3356–3369, Association for Computational Linguistics, 2020.

[8] European Union Commission, "Code of Conduct on Countering Illegal Hate Speech Online."
`https://ec.europa.eu/newsroom/just/item-detail.cfm?item_id=54300`,
2016.

[9] Twitter, "Hateful conduct policy." `https: //help.twitter.com/en/rules-and-policies/hateful-conduct-policy`.

[10] T. Davidson, D. Warmsley, M. W. Macy, and I. Weber, *Automated hate speech detection and the problem of offensive language*, in *Proceedings of the Eleventh International Conference on Web and Social Media, ICWSM 2017, Montréal, Québec, Canada, May 15-18, 2017*, pp. 512–515, AAAI Press, 2017.

[11] S. Walker, *Hate speech: The history of an American controversy.* U of Nebraska Press, 1994.

[12] J. B. Jacobs and K. Potter, *Hate crimes: Criminal law & identity politics.* Oxford University Press on Demand, 1998.

[13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, *Attention is all you need*, in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA* (I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, eds.), pp. 5998–6008, 2017.

[14] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, *et. al.*, *Language models are unsupervised multitask learners*, *OpenAI blog* **1** (2019), no. 8 9.

[15] J. Devlin, M. Chang, K. Lee, and K. Toutanova, *BERT: pre-training of deep bidirectional transformers for language understanding*, in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)* (J. Burstein, C. Doran, and T. Solorio, eds.), pp. 4171–4186, Association for Computational Linguistics, 2019.

[16] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter,

C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, *Language models are few-shot learners*, in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual* (H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, eds.), 2020.

[17] J. Qian, M. ElSherief, E. Belding, and W. Y. Wang, *Leveraging intra-user and inter-user representation learning for automated hate speech detection*, in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 118–123, 2018.

[18] Z. Waseem and D. Hovy, *Hateful symbols or hateful people? predictive features for hate speech detection on twitter*, in *Proceedings of the Student Research Workshop, SRW@HLT-NAACL 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pp. 88–93, The Association for Computational Linguistics, 2016.

[19] P. Badjatiya, S. Gupta, M. Gupta, and V. Varma, *Deep learning for hate speech detection in tweets*, in *Proceedings of the 26th International Conference on World Wide Web Companion*, pp. 759–760, International World Wide Web Conferences Steering Committee, 2017.

[20] J. Qian, M. ElSherief, E. M. Belding, and W. Y. Wang, *Hierarchical CVAE for fine-grained hate speech classification*, in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018* (E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, eds.), pp. 3550–3559, Association for Computational Linguistics, 2018.

[21] J. Qian, H. Wang, M. ElSherief, and X. Yan, *Lifelong learning of hate speech classification on social media*, in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021* (K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tür, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, eds.), pp. 2304–2314, Association for Computational Linguistics, 2021.

[22] J. Qian, M. ElSherief, E. M. Belding, and W. Y. Wang, *Learning to decipher hate symbols*, in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)* (J. Burstein, C. Doran, and T. Solorio, eds.), pp. 3006–3015, Association for Computational Linguistics, 2019.

[23] J. Qian, A. Bethke, Y. Liu, E. M. Belding, and W. Y. Wang, *A benchmark dataset for learning to intervene in online hate speech*, in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019* (K. Inui, J. Jiang, V. Ng, and X. Wan, eds.), pp. 4754–4763, Association for Computational Linguistics, 2019.

[24] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh, *Universal adversarial triggers for attacking and analyzing nlp*, in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2153–2162, 2019.

[25] S. Dathathri, A. Madotto, J. Lan, J. Hung, E. Frank, P. Molino, J. Yosinski, and R. Liu, *Plug and play language models: A simple approach to controlled text generation*, in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, OpenReview.net, 2020.

[26] B. Krause, A. D. Gotmare, B. McCann, N. S. Keskar, S. R. Joty, R. Socher, and N. F. Rajani, *Gedi: Generative discriminator guided sequence generation*, in *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021* (M. Moens, X. Huang, L. Specia, and S. W. Yih, eds.), pp. 4929–4952, Association for Computational Linguistics, 2021.

[27] J. Qian, L. Dong, Y. Shen, F. Wei, and W. Chen, *Controllable natural language generation with contrastive prefixes*, in *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022* (S. Muresan, P. Nakov, and A. Villavicencio, eds.), pp. 2912–2924, Association for Computational Linguistics, 2022.

[28] S. Hochreiter and J. Schmidhuber, *Long short-term memory*, *Neural computation* **9** (1997), no. 8 1735–1780.

[29] P. Burnap and M. L. Williams, *Us and Them: Identifying Cyber Hate on Twitter across Multiple Protected Characteristics*, *EPJ Data Science* **5** (2016), no. 1 11.

[30] A. Schmidt and M. Wiegand, *A Survey on Hate Speech Detection using Natural Language Processing*, in *SocialNLP'17: Proceedings of the 5th International Workshop on Natural Language Processing for Social Media*, pp. 1–10, 2017.

[31] C. Van Hee, E. Lefever, B. Verhoeven, J. Mennes, B. Desmet, G. De Pauw, W. Daelemans, and V. Hoste, *Detection and Fine-grained Classification of*

*Cyberbullying Events*, in *RANLP'15: International Conference Recent Advances in Natural Language Processing*, pp. 672–680, 2015.

[32] K. Dinakar, B. Jones, C. Havasi, H. Lieberman, and R. Picard, *Common Sense Reasoning for Detection, Prevention, and Mitigation of Cyberbullying*, *ACM Transactions on Interactive Intelligent Systems (TiiS)* **2** (2012), no. 3 18.

[33] G. Xiang, B. Fan, L. Wang, J. Hong, and C. Rose, *Detecting Offensive Tweets via Topical Feature Discovery over a Large Scale Twitter Corpus*, in *CIKM'12: Proceedings of the 21st ACM International Conference on Information and Knowledge Management*, pp. 1980–1984, ACM, 2012.

[34] M. Dadvar, D. Trieschnigg, R. Ordelman, and F. de Jong, *Improving cyberbullying detection with user context*, in *ECIR*, pp. 693–696, Springer, 2013.

[35] P. Zhou, W. Shi, J. Tian, Z. Qi, B. Li, H. Hao, and B. Xu, *Attention-based bidirectional long short-term memory networks for relation classification*, in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Short Papers)*, vol. 2, pp. 207–212, 2016.

[36] Y. Liu, C. Sun, L. Lin, and X. Wang, *Learning natural language inference using bidirectional LSTM model and inner-attention*, *CoRR* **abs/1605.09090** (2016) [arXiv:1605.0909].

[37] P. Indyk and R. Motwani, *Approximate nearest neighbors: towards removing the curse of dimensionality*, in *Proceedings of the 13th Annual ACM Symposium on Theory of Computing*, pp. 604–613, ACM, 1998.

[38] A. Gionis, P. Indyk, R. Motwani, *et. al.*, *Similarity search in high dimensions via hashing*, in *VLDB*, vol. 99, pp. 518–529, 1999.

[39] R. J. Williams, *Simple statistical gradient-following algorithms for connectionist reinforcement learning*, *Machine learning* **8** (1992), no. 3-4 229–256.

[40] J. A. Suykens and J. Vandewalle, *Least squares support vector machine classifiers*, *Neural processing letters* **9** (1999), no. 3 293–300.

[41] Y. Kim, *Convolutional neural networks for sentence classification*, in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL* (A. Moschitti, B. Pang, and W. Daelemans, eds.), pp. 1746–1751, ACL, 2014.

[42] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, *Autoencoding beyond pixels using a learned similarity metric*, in *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*, pp. 1558–1566, JMLR. org, 2016.

[43] K. Sohn, H. Lee, and X. Yan, *Learning structured output representation using deep conditional generative models*, in *Advances in Neural Information Processing Systems*, 2015.

[44] X. Zhang, J. Zhao, and Y. LeCun, *Character-level convolutional networks for text classification*, in *Advances in neural information processing systems*, pp. 649–657, 2015.

[45] R. Johnson and T. Zhang, *Effective use of word order for text categorization with convolutional neural networks*, in *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 103–112, 2015.

[46] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, *Recursive deep models for semantic compositionality over a sentiment treebank*, in *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631–1642, 2013.

[47] K. S. Tai, R. Socher, and C. D. Manning, *Improved semantic representations from tree-structured long short-term memory networks*, in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, vol. 1, pp. 1556–1566, 2015.

[48] D. Yogatama, C. Dyer, W. Ling, and P. Blunsom, *Generative and discriminative text classification with recurrent neural networks*, arXiv preprint arXiv:1703.01898 (2017).

[49] S. Lai, L. Xu, K. Liu, and J. Zhao, *Recurrent convolutional neural networks for text classification*, in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA* (B. Bonet and S. Koenig, eds.), pp. 2267–2273, AAAI Press, 2015.

[50] C. Zhou, C. Sun, Z. Liu, and F. Lau, *A c-lstm neural network for text classification*, arXiv preprint arXiv:1511.08630 (2015).

[51] D. Tang, B. Qin, and T. Liu, *Document modeling with gated recurrent neural network for sentiment classification*, in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015* (L. Màrquez, C. Callison-Burch, J. Su, D. Pighin, and Y. Marton, eds.), pp. 1422–1432, The Association for Computational Linguistics, 2015.

[52] Z. Yang, D. Yang, C. Dyer, X. He, A. J. Smola, and E. H. Hovy, *Hierarchical attention networks for document classification*, in *NAACL HLT 2016, The 2016*

*Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016* (K. Knight, A. Nenkova, and O. Rambow, eds.), pp. 1480–1489, The Association for Computational Linguistics, 2016.

[53] Q. Diao, M. Qiu, C.-Y. Wu, A. J. Smola, J. Jiang, and C. Wang, *Jointly modeling aspects, ratings and sentiments for movie recommendation (jmars)*, in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 193–202, ACM, 2014.

[54] D. P. Kingma and M. Welling, *Auto-encoding variational bayes*, in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.), 2014.

[55] T. Salimans, D. P. Kingma, and M. Welling, *Markov chain monte carlo and variational inference: Bridging the gap*, in *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015* (F. R. Bach and D. M. Blei, eds.), vol. 37 of *JMLR Workshop and Conference Proceedings*, pp. 1218–1226, JMLR.org, 2015.

[56] D. J. Rezende, S. Mohamed, and D. Wierstra, *Stochastic backpropagation and approximate inference in deep generative models*, in *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, vol. 32 of *JMLR Workshop and Conference Proceedings*, pp. 1278–1286, JMLR.org, 2014.

[57] B. Zhang, D. Xiong, J. Su, H. Duan, and M. Zhang, *Variational neural machine translation*, in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016* (J. Su, X. Carreras, and K. Duh, eds.), pp. 521–530, The Association for Computational Linguistics, 2016.

[58] T. Zhao, R. Zhao, and M. Eskénazi, *Learning discourse-level diversity for neural dialog models using conditional variational autoencoders*, in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers* (R. Barzilay and M. Kan, eds.), pp. 654–664, Association for Computational Linguistics, 2017.

[59] Southern Poverty Law Center, "Ideologies."
https://www.splcenter.org/fighting-hate/extremist-files/ideology.

[60] Southern Poverty Law Center, "Groups."
https://www.splcenter.org/fighting-hate/extremist-files/groups.

[61] L1ght, "Rising Levels of Hate Speech  Online Toxicity During This Time of Crisis."
https://l1ght.com/Toxicity_during_coronavirus_Report-L1ght.pdf, 2020.

[62] M. McCloskey and N. J. Cohen, *Catastrophic interference in connectionist networks: The sequential learning problem*, in *Psychology of learning and motivation*, vol. 24, pp. 109–165. Elsevier, 1989.

[63] R. Ratcliff, *Connectionist models of recognition memory: constraints imposed by learning and forgetting functions.*, *Psychological review* **97** (1990), no. 2 285.

[64] J. L. McClelland, B. L. McNaughton, and R. C. O'Reilly, *Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory.*, *Psychological review* **102** (1995), no. 3 419.

[65] R. M. French, *Catastrophic forgetting in connectionist networks*, *Trends in cognitive sciences* **3** (1999), no. 4 128–135.

[66] D. M. Blei, A. Y. Ng, and M. I. Jordan, *Latent dirichlet allocation*, *Journal of machine Learning research* **3** (2003) 993–1022.

[67] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, *et. al.*, *Overcoming catastrophic forgetting in neural networks*, *Proceedings of the national academy of sciences* **114** (2017), no. 13 3521–3526.

[68] F. Zenke, B. Poole, and S. Ganguli, *Continual learning through synaptic intelligence*, in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017* (D. Precup and Y. W. Teh, eds.), vol. 70 of *Proceedings of Machine Learning Research*, pp. 3987–3995, 2017.

[69] C. Fernando, D. Banarse, C. Blundell, Y. Zwols, D. Ha, A. A. Rusu, A. Pritzel, and D. Wierstra, *Pathnet: Evolution channels gradient descent in super neural networks*, *CoRR* **abs/1701.08734** (2017).

[70] X. Liu, M. Masana, L. Herranz, J. Van de Weijer, A. M. Lopez, and A. D. Bagdanov, *Rotate your networks: Better weight consolidation and less catastrophic forgetting*, in *2018 24th International Conference on Pattern Recognition (ICPR)*, pp. 2262–2268, IEEE, 2018.

[71] H. Ritter, A. Botev, and D. Barber, *Online structured laplace approximations for overcoming catastrophic forgetting*, in *Advances in Neural Information Processing Systems*, pp. 3738–3748, 2018.

[72] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, *Progressive neural networks*, *CoRR* **abs/1606.04671** (2016).

[73] J. L. Part and O. Lemon, *Incremental on-line learning of object classes using a combination of self-organizing incremental neural networks and deep convolutional neural networks*, in *Workshop on Bio-inspired Social Robot Learning in Home Scenarios (IROS), Daejeon, Korea*, 2016.

[74] J. L. Part and O. Lemon, *Incremental online learning of objects for robots operating in real environments*, in *2017 Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*, pp. 304–310, IEEE, 2017.

[75] D. Lopez-Paz and M. Ranzato, *Gradient episodic memory for continual learning*, in *Advances in Neural Information Processing Systems*, pp. 6467–6476, 2017.

[76] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny, *Efficient lifelong learning with A-GEM*, in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, OpenReview.net, 2019.

[77] C. de Masson d'Autume, S. Ruder, L. Kong, and D. Yogatama, *Episodic memory in lifelong language learning*, in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada* (H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, eds.), pp. 13122–13131, 2019.

[78] F. Sun, C. Ho, and H. Lee, *LAMOL: language modeling for lifelong language learning*, in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.

[79] R. Xia, J. Jiang, and H. He, *Distantly supervised lifelong learning for large-scale social media sentiment analysis*, *IEEE Trans. Affect. Comput.* **8** (2017), no. 4 480–491.

[80] H. Wang, W. Xiong, M. Yu, X. Guo, S. Chang, and W. Y. Wang, *Sentence embedding alignment for lifelong relation extraction*, in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 796–806, 2019.

[81] J. Golbeck, Z. Ashktorab, R. O. Banjo, A. Berlinger, S. Bhagwan, C. Buntain, P. Cheakalos, A. A. Geller, Q. Gergory, R. K. Gnanasekaran, *et. al.*, *A large*

*labeled corpus for online harassment research*, in *Proceedings of the 2017 ACM on Web Science Conference*, pp. 229–233, ACM, 2017.

[82] H. Zhang, X. Xiao, and O. Hasegawa, *A load-balancing self-organizing incremental neural network*, IEEE Transactions on Neural Networks and Learning Systems **25** (2013), no. 6 1096–1105.

[83] K. Cho, B. van Merrienboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, *Learning phrase representations using RNN encoder-decoder for statistical machine translation*, in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pp. 1724–1734, 2014.

[84] I. Sutskever, O. Vinyals, and Q. V. Le, *Sequence to sequence learning with neural networks*, in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada* (Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, eds.), pp. 3104–3112, 2014.

[85] B. Han and T. Baldwin, *Lexical normalisation of short text messages: Makn sens a# twitter*, in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pp. 368–378, Association for Computational Linguistics, 2011.

[86] B. Han, P. Cook, and T. Baldwin, *Automatically constructing a normalisation dictionary for microblogs*, in *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pp. 421–432, Association for Computational Linguistics, 2012.

[87] F. Liu, F. Weng, and X. Jiang, *A broad-coverage normalization system for social media language*, in *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pp. 1035–1044, Association for Computational Linguistics, 2012.

[88] G. Chrupała, *Normalizing tweets with edit scripts and recurrent neural embeddings*, in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, vol. 2, pp. 680–686, 2014.

[89] K. Knight, A. Nair, N. Rathod, and K. Yamada, *Unsupervised analysis for decipherment problems*, in *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pp. 499–506, 2006.

[90] S. Ravi and K. Knight, *Deciphering foreign language*, in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pp. 12–21, Association for Computational Linguistics, 2011.

[91] D. Bahdanau, K. Cho, and Y. Bengio, *Neural machine translation by jointly learning to align and translate*, in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.), 2015.

[92] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, *et. al.*, *Google's neural machine translation system: Bridging the gap between human and machine translation*, arXiv preprint arXiv:1609.08144 (2016).

[93] K. Ni and W. Y. Wang, *Learning to explain non-standard english words and phrases*, in *Proceedings of the Eighth International Joint Conference on Natural Language Processing, IJCNLP 2017, Taipei, Taiwan, November 27 - December 1, 2017, Volume 2: Short Papers* (G. Kondrak and T. Watanabe, eds.), pp. 413–417, Asian Federation of Natural Language Processing, 2017.

[94] F. Hill, K. Cho, A. Korhonen, and Y. Bengio, *Learning to understand phrases by embedding the dictionary*, *Transactions of the Association of Computational Linguistics* **4** (2016) 17–30.

[95] T. Noraset, C. Liang, L. Birnbaum, and D. Downey, *Definition modeling: Learning to define word embeddings in natural language*, in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA* (S. Singh and S. Markovitch, eds.), pp. 3259–3266, AAAI Press, 2017.

[96] J. Pavlopoulos, P. Malakasiotis, and I. Androutsopoulos, *Deeper attention to abusive user content moderation*, in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 1125–1135, 2017.

[97] L. Gao, A. Kuppersmith, and R. Huang, *Recognizing explicit and implicit hate speech using a weakly supervised two-path bootstrapping approach*, in *Proceedings of the Eighth International Joint Conference on Natural Language Processing*, pp. 774–782, 2017.

[98] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*, in *Proceedings of the 3rd International Conference on Learning Representations*, 2015.

[99] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, *BLEU: a method for automatic evaluation of machine translation*, in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pp. 311–318, Association for Computational Linguistics, 2002.

[100] C.-Y. Lin, *Rouge: A package for automatic evaluation of summaries*, *Text Summarization Branches Out* (2004).

[101] S. Banerjee and A. Lavie, *Meteor: An automatic metric for MT evaluation with improved correlation with human judgments*, in *Proceedings of the ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.

[102] M. ElSherief, V. Kulkarni, D. Nguyen, W. Y. Wang, and E. Belding, *Hate lingo: A target-based linguistic analysis of hate speech in social media*, in *Twelfth International AAAI Conference on Web and Social Media*, 2018.

[103] M. ElSherief, S. Nilizadeh, D. Nguyen, G. Vigna, and E. Belding, *Peer to peer hate: Hate speech instigators and their targets*, in *Twelfth International AAAI Conference on Web and Social Media*, 2018.

[104] A. Ohlheiser, *Banned from twitter? this site promises you can say whatever you want*, *The Washington Post* (2016).

[105] Pew Research Center, "Online Harassment 2017." `https://www.pewresearch.org/internet/2017/07/11/online-harassment-2017/`.

[106] H. Bielefeldt, F. La Rue, and G. Muigai, *Ohchr expert workshops on the prohibition of incitement to national, racial or religious hatred*, *Expert workshop on the Americas* (2011).

[107] D. Chatzakou, N. Kourtellis, J. Blackburn, E. De Cristofaro, G. Stringhini, and A. Vakali, *Mean birds: Detecting aggression and bullying on twitter*, in *Proceedings of the 2017 ACM on web science conference*, pp. 13–22, ACM, 2017.

[108] G. W. Kennedy III, A. W. McCollough, E. Dixon, A. Bastidas, J. Ryan, C. Loo, and S. Sahay, *Technology solutions to combat online harassment*, in *Proceedings of the first workshop on abusive language online*, pp. 73–77, 2017.

[109] A. M. Founta, C. Djouvas, D. Chatzakou, I. Leontiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kourtellis, *Large scale crowdsourcing and characterization of twitter abusive behavior*, in *Twelfth International AAAI Conference on Web and Social Media*, 2018.

[110] VICE, "Here Are Reddit's Whiniest, Most Low-Key Toxic Subreddits." `https://www.vice.com/en_us/article/8xxymb/here-are-reddits-whiniest-most-low-key-toxic-subreddits`, 2017.

[111] Facebook, "Hate Speech."
https://m.facebook.com/communitystandards/hate_speech/, 2019.

[112] J. Li, W. Monroe, A. Ritter, D. Jurafsky, M. Galley, and J. Gao, *Deep reinforcement learning for dialogue generation*, in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pp. 1192–1202, 2016.

[113] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, *Distributed representations of words and phrases and their compositionality*, in *Advances in neural information processing systems*, pp. 3111–3119, 2013.

[114] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Józefowicz, and S. Bengio, *Generating sentences from a continuous space*, in *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, pp. 10–21, 2016.

[115] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba, *Sequence level training with recurrent neural networks*, in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.), 2016.

[116] R. Paulus, C. Xiong, and R. Socher, *A deep reinforced model for abstractive summarization*, in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, OpenReview.net, 2018.

[117] X. L. Li and P. Liang, *Prefix-tuning: Optimizing continuous prompts for generation*, in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021* (C. Zong, F. Xia, W. Li, and R. Navigli, eds.), pp. 4582–4597, Association for Computational Linguistics, 2021.

[118] N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher, *CTRL: A conditional transformer language model for controllable generation*, *CoRR* **abs/1909.05858** (2019) [arXiv:1909.0585].

[119] Z. Hu, Z. Yang, X. Liang, R. Salakhutdinov, and E. P. Xing, *Toward controlled generation of text*, in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017* (D. Precup and Y. W. Teh, eds.), vol. 70 of *Proceedings of Machine Learning Research*, pp. 1587–1596, PMLR, 2017.

[120] J. Ficler and Y. Goldberg, *Controlling linguistic style aspects in neural language generation*, *CoRR* **abs/1707.02633** (2017) [arXiv:1707.0263].

[121] A. Holtzman, J. Buys, M. Forbes, A. Bosselut, D. Golub, and Y. Choi, *Learning to write with cooperative discriminators*, in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers* (I. Gurevych and Y. Miyao, eds.), pp. 1638–1649, Association for Computational Linguistics, 2018.

[122] K. Yang and D. Klein, *FUDGE: controlled text generation with future discriminators*, in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021* (K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tür, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, eds.), pp. 3511–3535, Association for Computational Linguistics, 2021.

[123] Z. Lin and M. O. Riedl, *Plug-and-blend: A framework for plug-and-play controllable story generation with sketches*, in *Proceedings of the Seventeenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment, AIIDE 2021, virtual, October 11-15, 2021* (D. Thue and S. G. Ware, eds.), pp. 58–65, AAAI Press, 2021.

[124] D. Yu, Z. Yu, and K. Sagae, *Attribute alignment: Controlling text generation from pre-trained language models*, in *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021* (M. Moens, X. Huang, L. Specia, and S. W. Yih, eds.), pp. 2251–2268, Association for Computational Linguistics, 2021.

[125] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, *Neural discrete representation learning*, in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA* (I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, eds.), pp. 6306–6315, 2017.

[126] E. Jang, S. Gu, and B. Poole, *Categorical reparameterization with gumbel-softmax*, in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, OpenReview.net, 2017.

[127] C. J. Maddison, A. Mnih, and Y. W. Teh, *The concrete distribution: A continuous relaxation of discrete random variables*, in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, OpenReview.net, 2017.

[128] C. K. Sønderby, B. Poole, and A. Mnih, *Continuous relaxation training of discrete latent variable image models*, in *Beysian DeepLearning workshop, NIPS*, vol. 201, 2017.

[129] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, *Zero-shot text-to-image generation*, in *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event* (M. Meila and T. Zhang, eds.), vol. 139 of *Proceedings of Machine Learning Research*, pp. 8821–8831, PMLR, 2021.

[130] Jigsaw, "Jigsaw Toxic Comment Classification Challenge Dataset." `https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/`.

[131] A. Baheti, M. Sap, A. Ritter, and M. Riedl, *Just say no: Analyzing the stance of neural dialogue generation in offensive contexts*, in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4846–4862, 2021.

[132] M. Zampieri, P. Nakov, S. Rosenthal, P. Atanasova, G. Karadzhov, H. Mubarak, L. Derczynski, Z. Pitenis, and Ç. Çöltekin, *Semeval-2020 task 12: Multilingual offensive language identification in social media (offenseval 2020)*, in *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pp. 1425–1447, 2020.

[133] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, *Roberta: A robustly optimized BERT pretraining approach*, *CoRR* **abs/1907.11692** (2019) [arXiv:1907.1169].

[134] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov, *Unsupervised cross-lingual representation learning at scale*, in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020* (D. Jurafsky, J. Chai, N. Schluter, and J. R. Tetreault, eds.), pp. 8440–8451, Association for Computational Linguistics, 2020.

[135] S. Wang, J. Liu, X. Ouyang, and Y. Sun, *Galileo at semeval-2020 task 12: Multi-lingual learning for offensive language identification using pre-trained language models*, in *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pp. 1448–1455, 2020.

[136] Q. Zhu, Z. Lin, Y. Zhang, J. Sun, X. Li, Q. Lin, Y. Dang, and R. Xu, *Hitsz-hlt at semeval-2021 task 5: Ensemble sequence labeling and span boundary detection for toxic span detection*, in *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pp. 521–526, 2021.

[137] V. A. Nguyen, T. M. Nguyen, H. Q. Dao, and Q. H. Pham, *S-nlp at semeval-2021 task 5: An analysis of dual networks for sequence tagging*, in *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pp. 888–897, 2021.

[138] L. Hanu and Unitary team, "Detoxify."
`https://github.com/unitaryai/detoxify`, 2020.

[139] B. Kennedy, X. Jin, A. M. Davani, M. Dehghani, and X. Ren, *Contextualizing hate speech classifiers with post-hoc explanation*, in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5435–5442, 2020.

[140] Y. Zhang, S. Sun, M. Galley, Y.-C. Chen, C. Brockett, X. Gao, J. Gao, J. Liu, and B. Dolan, *Dialogpt: Large-scale generative pre-training for conversational response generation*, in *ACL, system demonstration*, 2020.

[141] Y. Li and C. Scarton, *Revisiting rumour stance classification: Dealing with imbalanced data*, in *Proceedings of the 3rd International Workshop on Rumours and Deception in Social Media (RDSM)*, pp. 38–44, 2020.

[142] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, *et. al.*, *Transformers: State-of-the-art natural language processing*, in *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.

[143] F. Yang, X. Peng, G. Ghosh, R. Shilon, H. Ma, E. Moore, and G. Predovic, *Exploring deep multimodal fusion of text and photo for hate speech classification*, in *Proceedings of the third workshop on abusive language online*, pp. 11–18, 2019.