

On the Importance and Challenges of the Experimental Design of Multilingual Toxic Content Detection”

by

Nedjma Djouhra OUSIDHOUM

A Thesis Submitted to
The Hong Kong University of Science and Technology
in Partial Fulfillment of the Requirements for
the Degree of Doctor of Philosophy
in

August 2021, Hong Kong

Authorization

I hereby declare that I am the sole author of the thesis.

I authorize the Hong Kong University of Science and Technology to lend this thesis to other institutions or individuals for the purpose of scholarly research.

I further authorize the Hong Kong University of Science and Technology to reproduce the thesis by photocopying or by other means, in total or in part, at the request of other institutions or individuals for the purpose of scholarly research.

Nedjma Djouhra OUSIDHOUM

August 2021

On the Importance and Challenges of the Experimental Design of Multilingual Toxic Content Detection”

by

Nedjma Djouhra OUSIDHOUM

This is to certify that I have examined the above PhD thesis
and have found that it is complete and satisfactory in all respects,
and that any and all revisions required by
the thesis examination committee have been made.

Dr. Yangqiu SONG, Thesis Supervisor
Department of Computer Science and Engineering

Prof. Dit-Yan YEUNG, Thesis Co-Supervisor
Department of Computer Science and Engineering

Prof. Dit-Yan YEUNG
Head, Department of Computer Science and Engineering

Department of Computer Science and Engineering, HKUST

August 2021

This thesis is for the marginalized, the harmed, the ignored, the abused, and the dehumanized.

*To me, you are more than statistics and case studies in reports. May we and the future
generations build a just, kind, and better world.*

*I also dedicate this thesis to my extraordinary mama Amina Rahal who taught me that a selfish
existence is a pointless one. Nul hommage ne sera à ta hauteur. Merci maman, je t'aime plus
que tout au monde... شكرا ماما العزيزة*

Acknowledgments

First and foremost, I would like to thank my supervisors Dr. Yangqiu Song for his guidance and remarks and Prof. Dit-Yan Yeung for helping me through the process since 2018.

I would like to thank my examiners Prof. Pascale Fung, Dr. Brian Mak, Prof. Nevin Zhang from HKUST, and Dr. Preslav Nakov from HBKU and QCRI for doing me the honor of being part of my PhD thesis committee and for their feedback and helpful remarks.

Thanks to my collaborators Tianqing Fang, Zizheng Lin, and Xinran Zhao and to my labmates for being part of my PhD journey. I learned from you during our meetings and reading groups.

I would like to thank my former colleagues and friends Dr. Meriem Beloucif, Dr. Markus Saers, and Dr. Chi-kiu (Jackie) Lo with whom I shared a unique experience! During the first four years of my PhD, Markus was a remarkable senior labmate from whom I learned a lot. Thank you, Markus! Also, thank you Jackie for your help with my professional life.

Throughout this journey, Dr. Meriem Beloucif was much more than a former labmate. Thank you Meriem for being a brilliant and supportive colleague, but more importantly an exceptional and dear friend who changed my black and white ideas on friendship. Thank you for your loyalty, for believing in me when I had a hard time doing it myself, for sticking around during some of the hardest moments in my life, and for being a shoulder to cry on despite the distance. We have been like sisters along this journey and I am proud to be your daughter's *tata* and to love her as much I would have loved my own niece.

A million thanks to my mentor Professor Nacéra Bensaou from USTHB who has been the one who changed my life's trajectory, inspired me, truly believed in me, and encouraged me to do research. She taught me resilience (among other things) and I hope that one day, I can be someone's true mentor the way she was mine. Besides that, I am proud to be her friend and I am grateful to benefit from her incredible knowledge and wisdom until now.

Many thanks to our kind communication tutor Ms. Shauna Dalton and Dr. Meriem Beloucif for proofreading my thesis. I would also like to thank Professor Brahim Bensaou for his professional advises, my brilliant (remote) teacher and friend Alaa Abuarab for her beautiful writings and classes, and my amazing counselor Ms. Vava Kwok who was incredibly helpful throughout this journey.

For more than five years, I have been volunteering as a TA in Po Leung Keung. I would like to thank each kid in my classes for being smart, lively, funny, adorable and for being the highlight of so many weeks. If any of you comes across my thesis one day or attends HKUST, I

would like to tell you that each one of you has had a special place in my heart, I still remember your names, and I still have the cards you drew for me in an envelope that I cherish.

HKUST would not have been the way it is without so many people behind the scenes. I would like to thank Ms. Connie Lau from the CSE admin office for being exceptionally helpful. I would also like to thank many cleaning ladies, security agents, and café barista for being extremely nice, warm, and making my days with their kindness especially during the pandemic.

Last but not least, a big thank you to my family back home and abroad. Thanks to my uncles and aunties for checking up on me and helping me before I even asked. Thanks to momani the loveliest grandma in the world, my supportive father Hocine, and my two brothers and cheerleaders Abdou and Raouf for their unconditional love and valuable support.

And thank you mama for being my wonderful mother and friend. This thesis is dedicated to you.

Contents

Title Page	i
Authorization Page	ii
Signature Page	iii
Acknowledgments	v
Table of Contents	vii
List of Figures	xi
List of Tables	xiv
Abstract	xviii
1 Introduction	1
1.1 Overview of Hate Speech Detection	1
1.1.1 Defining Hate Speech	2
1.1.2 Online Policy Against Hate Speech	2
1.2 The Importance of Multilingual Hate Speech Detection	3
1.3 Hate Speech Annotations	4
1.4 A Cultural Study on Hate Speech	4
1.5 Selection Bias in Hate Speech	6
1.6 Probing Toxic Content in Large Pre-trained Language Models	7
1.7 Thesis Organization	9
2 Background	10
2.1 Automatic Detection of Hate Speech and Toxic Language	10
2.1.1 Hate Speech and Toxic Language	10

2.1.2	Defining Hate Speech for Automatic Detection	10
2.2	Language Resources	11
2.2.1	Lexicons	11
2.2.2	Datasets	11
2.3	Toxic Language Classification	12
2.3.1	Coarse-Grained Toxic Language Classification	13
2.3.2	Fine-Grained Toxic Language Classification	13
2.3.3	Word Representations for Classification	14
2.4	Bias in NLP	14
2.4.1	Selection Bias	15
2.4.2	Label Bias	15
2.4.3	Semantic Bias	15
2.4.4	Model Overamplification	16
2.4.5	Solutions to Bias in Social Data	16
2.5	Toxic Language Classifiers	16
2.5.1	Baselines	17
2.5.2	Deep Learning Models	17
3	Multilingual Multi-Aspect Hate Speech Detection	18
3.1	Dataset Construction	18
3.1.1	Data Collection	19
3.1.2	Annotation Challenges	20
3.1.3	Annotation Process	21
3.1.4	Pilot Dataset	22
3.1.5	Final Dataset	22
3.2	Experiments	25
3.2.1	Models	26
3.2.2	Results and Analysis	26
4	Cultural Differences in Hate Speech	30
4.1	Cultural Studies in NLP	31
4.2	Overview of the Data	32
4.2.1	Description of the Datasets	32

4.2.2	Hate Speech Aspects	32
4.3	Analysis and Discussion	34
4.3.1	Frequent Words in Hateful Tweets	34
4.3.2	Comparison Between Datasets	37
4.3.3	Coherence Scores	38
4.3.4	Case Study on Annotators' Reactions	39
4.3.5	Case Study on Similar Targets	40
5	Selection Bias in Hate Speech Detection	42
5.1	Topic Modeling	42
5.2	Bias Estimation	43
5.2.1	Predefined Keywords	43
5.2.2	Topic Models	44
5.2.3	Bias Metrics	46
5.3	Results	48
5.3.1	Experimental Settings	48
5.3.2	Robustness Towards The Variability of Topic Distribution	48
5.3.3	Robustness of Keyword-based Selection	49
5.3.4	Hate Speech Embeddings	49
5.3.5	General versus Corpus-Specific Lists of Keywords	51
5.3.6	WordNet and Targeted Hate Bias	51
5.3.7	Case Study	52
5.4	Discussion	53
6	Probing Toxic Content in Large Pre-Trained Language Models	55
6.1	Methodology	56
6.1.1	Probing Patterns	56
6.1.2	Lists of Social Groups	57
6.1.3	The Generated Data	58
6.1.4	Probing Classifiers	59
6.1.5	Bias in Toxic Language Classifiers	59
6.2	Experiments	60
6.2.1	Main Results	60

6.2.2	Human Evaluation	62
6.3	A Case Study On Offensive Content Generated by PTLMs	63
6.4	Frequent Content Analysis	64
6.4.1	Frequent Content in English	64
6.4.2	Frequent Content in French and Arabic	65
6.4.3	Ethical Considerations	65
7	Conclusion	69
	References	71
	Publications	82

List of Figures

1.1	Examples of tweets where (1) immigrants are accused of harming society without the use of any direct insult; (2) a Hispanic person is insulted using a slur; and (3) a slur is used to give a personal account. The three examples show the complexity of hate speech and prove that profanity is not a clear indicator of hate speech.	2
1.2	Arabic and French tweets targeting the same group of people using a different vocabulary. The French and Arabic tweets demonstrate the complexity of the task and draw the distinction between hate speech directed to the same community, in the two languages, due to different socio-cultural backgrounds.	3
1.3	Annotated English example with respect to different aspects, namely, directness, hostility, target attribute, target group, and annotator’s sentiment.	4
2.1	Sources of Bias in NLP applications. While Shah et al. (2020) present a thorough theoretical framework, we show general sources of bias in NLP models which are commonly noticed in hate speech and toxic detection.	15
3.1	Multi-aspect annotations in our dataset. We show different annotated multi-labeled aspects.	19
4.1	Three hateful examples in English, French and Arabic, targeting immigrants by describing them as traitors in a nationalist tweet in Arabic, invaders in a French tweet, and people who carry diseases in English.	31
4.2	Top words in hateful tweets in three English datasets.	35
4.3	Top words in hateful tweets in three different Arabic datasets.	36
4.4	Frequent words in Ibrohim and Budi (2019)’s Indonesian hateful tweets.	36
4.5	Top words in some non-English hateful tweets in languages largely spoken in the EU.	37

- 4.6 coherence score variations for different datasets when generating 8 topics containing words in the interval [2, 100]. **AR1**, **AR2**, **AR3** refer to the Arabic datasets by Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019) respectively, **DE** to the German dataset by Ross et al. (2017), **EN1**, **EN2**, **EN3** to the English datasets by Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016) respectively, **FR** to the French dataset by Ousidhoum et al. (2019), **IT** to the Italian dataset by Sanguinetti et al. (2018), **ID** to the Indonesian dataset by Ibrohim and Budi (2019), and **PT** to the Portuguese dataset by Fortuna et al. (2019). 39
- 5.1 Average B_1 scores based on topic and word numbers in the interval [2, 100]. We fix the number of topics to 8 when we alter the number of words and similarly, we fix the number of words to 8 when we change the number of topics. We use the **multilingual Babylon embeddings** to compute the semantic similarity between words. 47
- 5.2 Average B_2 scores based on topic and word numbers in the interval [2, 100]. We fix the number of topics to 8 when we alter the number of words and similarly, we fix the number of words to 8 when we change the number of topics. We use the **multilingual Babylon embeddings** to compute the semantic similarity between words. 47
- 5.3 B_1 score variations for different datasets. The numbers of topics and words in topics are in the range [2, 100]. We use **multilingual Babylon embeddings** to compute the semantic similarity between words. EN1, EN2, EN3 refer to Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016); and AR1, AR2, AR3 to Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019), respectively. 49
- 5.4 B_2 score variations for different datasets. The numbers of topics and words in topics are in the range [2, 100]. We use **multilingual Babylon embeddings** to compute the semantic similarity between words. EN1, EN2, EN3 refer to Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016); and AR1, AR2, AR3 to Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019), respectively. 50

5.5	Variations of B_1 (in blue) and B_2 (in red) scores on the German and Indonesian datasets.	53
6.1	An example of generated content using BERT. Intuitively, one would think that adjectives would be prioritized over ethnic/religious affiliations in a cause-effect cloze statement, which appears not to be the case. Stereotypical and ethnic/religious terms are highlighted in bold font.	56

List of Tables

1.1	Examples of keywords present in the predefined lists along with their English translations. The keywords include terms frequently associated with controversies such as <i>communist</i> in Italian, slurs such as <i>m*ng*</i> in French, insults such as <i>pig</i> in Arabic, and hashtags such as <i>rapefugees</i> in German.	6
1.2	Confusing examples where the last word is predicted by a PTLM.	8
1.3	Insulting examples, where the last word is predicted by a PTLM. Sentences include offensive content, implicit insults, microaggressions, and stereotypes.	8
2.1	Description of different available hate speech and offensive language datasets.	12
2.2	Examples of toxic language datasets with coarse-grained labels.	13
2.3	Examples of fine-grained labeling schemes with numbers of labels per annotated aspect.	13
3.1	The label distributions of each task. The counts of direct and indirect hate speech include all tweets except those that are single-labeled as “normal”. Hostility type and annotator’s sentiment are multilabel classification tasks, while target attribute and target group are not. We show the counts of the top 5 target groups among 16 in total.	23
3.2	Full evaluation scores of the only binary classification task where the single task single language model consistently outperforms multilingual multitask models.	27
3.3	Full evaluation of tasks where multilingual and multitask models outperform on average single-task-single-language model on four different tasks.	28

4.1	The number of collected tweets (#POSTS), the size of vocabulary (Vocabulary), the average size of tweets (Tweet) in eleven datasets. AR1, AR2, AR3 refer to the Arabic datasets by Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019) respectively, DE to the German dataset by Ross et al. (2017), EN1, EN2, EN3 to the English datasets by Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016) respectively, FR to the French dataset by Ousidhoum et al. (2019), IT to the Italian dataset by Sanguinetti et al. (2018), ID to the Indonesian dataset by Ibrohim and Budi (2019), and PT to the Portuguese dataset by Fortuna et al. (2019).	33
4.2	Examples of 5 word topics generated by LDA from each of the previously described datasets.	38
4.3	Top 10 word topics generated by LDA based on the annotators' reactions to hateful tweets. Arabic words are written from right to left. Hence, the translations are shown in the reverse order.	40
4.4	Topic words in two Arabic datasets that discriminate people based on religious affiliations. Arabic words are written from right to left, therefore the translations are shown in the reverse order.	41
5.1	Examples of keywords present in the predefined lists of keywords and their English translations. The keywords include terms frequently associated with controversies, demeaning terms, and hashtags.	44
5.2	Examples of topics of length 3 generated by LDA. Non-English topics are presented with their English translations. Some topics contain slurs, named entities, and hashtags.	45
5.3	B₁ scores based on trained hate speech embeddings for 10 topics. We have manually clustered the keywords released with our dataset Ousidhoum et al. (2019) based on discriminating target attributes. For instance, the word <i>ni**er</i> belongs the ORIGIN category, <i>raghead</i> to RELIGION , and <i>c**t</i> to GENDER . For normalization purposes, we skipped disability since we did not find predefined Arabic keywords that target people with disabilities.	50

- 5.4 B_1 scores for English hate speech datasets using **WordNet** given 10 topics and keywords clustered based on **ORIGIN**, **RELIGION**, and **GENDER**. The scores are reported for tweets that have not been labeled *non-hateful* or *normal*. Although we initially attempted to study the differences of pre-trained word embeddings and word associations, we found that many (w_j, w'_k) pairs involve out-of-the-vocabulary words. In such cases, $Sim(w_j, w'_k)$ would have a WordNet similarity score $WUP = 0$ which is why the scores are in the range $[0.25, 0.35]$. 52
- 5.5 Given the average B_1 and B_2 scores generated for each dataset, based on topics (**#TOPICS**) and topic words (**#WORDS**) in the interval $[2, 100]$, respectively, we compute Spearman’s correlation scores between B_1 and B_2 and (1) the number of keywords $|w'|$ and average cosine similarity between keywords w'_{sim} given the language of the dataset; in addition to (2) the number of collected tweets **#TWEETS**, their average size **TWEET**, and size of vocabulary **VOCAB** in each dataset. 54
- 6.1 Patterns used with the **ATOMIC** actions. Given the nature of PTLMs and for the sake of our multilingual study, we use the pronouns *he* and *she* even for *PersonX*. *ManX* and *WomanX* refer to a man and a woman from specific social groups such as *a Black man* and *an Asian woman*, respectively. 56
- 6.2 Examples of social groups we use in our experiments. **Race** refers to different racial groups; **Religion** to different (non)religious affiliations; **Gender** to different genders and sexual orientations; **Politics** to various political views; **Intersectional** to social groups that fall into the intersection of two attributes such as gender and race; and **Marginalized** to commonly marginalized communities. 57
- 6.3 Examples of top 10 predicted reasons given various social groups and actions. 58
- 6.4 F1 and Accuracy scores of the logistic regression (LR) toxic language classifiers. 59

- 6.5 Proportions of the generated sentences which are classified as *toxic* by the LR classifiers. $\%@k$ refers to the proportion of toxic sentences when retrieving top k words predicted by the corresponding PTLM. BERT tends to generate more potentially toxic content compared to GPT-2 and RoBERTa, which may be due to the fact that GPT-2 generates a large number of stop words and punctuation marks. The variations across languages are largely due to the difference in the sizes of the evaluation samples, since we have fewer instances to assess in French and Arabic. In addition, the French classifier is trained on only one relatively small dataset. 60
- 6.6 The scores in this table indicate the proportions of potentially toxic statements with respect to a given social group based on content generated by different PTLMs. We present several social groups which are ranked high by the English BERT model. 61
- 6.7 Human Evaluation of 100 predicted sentences by BERT, CamemBERT, and AraBERT labeled by five annotators. **#Insult** refers to problematic examples considered as insulting, **#Stereotype** refers to stereotypical content, **#Confusing** to confusing content and **#Normal** to normal content. The Fleiss Kappa scores are 0.63 for English, 0.64 for French, and 0.21 for Arabic. 62
- 6.8 Frequency (**F**) of Social groups (**S**) associated with names of animals in the predictions. The words are sometimes brought up as a reason (**e.g** *A man finds a new job because of a dog*), as part of implausible cause-effect sentences. Yet, sometimes they are used as direct insults (**e.g** *because he is a dog*). The last statement is insulting in Arabic. 64
- 6.9 Examples of relatively informative descriptive nouns and adjectives which appear as Top-1 predictions. We show the two main social groups that are associated with them. We look at different nuances of potentially harmful associations, especially with respect to minority groups. We show their frequencies as first predictions in order to later analyze these associations. 67
- 6.10 Arabic and French examples of relatively informative noun and adjective Top-1 predictions within the two main social groups which are associated with them. 68

On the Importance and Challenges of the Experimental Design of Multilingual Toxic Content Detection”

by Nedjma Djouhra OUSIDHOUM

Department of Computer Science and Engineering, HKUST

The Hong Kong University of Science and Technology

Abstract

With the expanding use of social media platforms such as Twitter and the amount of text data generated online, hate speech and toxic language have been proven to negatively affect individuals in general, and marginalized communities in particular. In order to improve the online moderation process, there has been an increasing need for accurate detection tools which do not only flag bad words but rather help to filter out toxic content in a more nuanced fashion. Hence, a problem of central importance is to acquire data of better quality in order to train toxic content detection models. However, the absence of a universal definition of hate speech makes the collection process hard and the training corpora sparse, imbalanced, and challenging for current machine learning techniques. In this thesis, we address the problem of automatic toxic content detection along three main axes: (1) the construction of resources lacking in robust toxic language and hate speech detection systems, (2) the study of bias in hate speech and toxic language classifiers, and (3) the assessment of inherent toxicity and harmful biases within NLP systems by looking into Large Pre-trained Language Models (PTLMs), which are at the core of these systems.

In order to train a multi-cultural, fine-grained hate speech and toxic content detection system, we have built a new multi-aspect hate speech dataset in English, French, and Arabic. We also provide a detailed annotation scheme, which indicates (a) whether a tweet is direct or indirect; (b) whether it is offensive, disrespectful, hateful, fearful out of ignorance, abusive, or normal; (c) the attribute based on which it discriminates against an individual or a group of people; (d) the name of this group; and (e) how annotators feel about this tweet given a range of negative to neutral sentiments. We define classification tasks based on each labeled aspect and use multi-task learning to investigate how such a paradigm can improve the detection process.

Unsurprisingly, when testing the detection system, the imbalanced data along with implicit toxic content and misleading instances has resulted in false positives and false negatives. We examine misclassification instances due to the frequently neglected yet deep-rooted *selection bias* caused by the data collection process. In contrast to work on bias, which typically focuses on the classification performance, we investigate another source of bias and present two language and label-agnostic evaluation metrics based on topic models and semantic similarity measures to evaluate the extent of such a problem on various datasets. Furthermore, since we generally focus on English and overlook other languages, we notice a gap in content moderation across languages and cultures, especially in low-resource settings. Hence, we leverage the observed differences and correlations across languages, datasets, and annotation schemes to carry a study on multilingual toxic language data and how people react to it.

Finally, social media posts are part of the training data of Large Pre-trained Language Models (PTLMs), which are at the center of all major NLP systems nowadays. Despite their incontestable usefulness and effectiveness, PTLMs have been shown to carry and reproduce harmful biases due to the sources of their training data among other reasons. We propose a methodology to probe the potentially toxic content that they convey with respect to a set of templates, and report how often they enable toxicity towards specific communities in English, French, and Arabic.

The results presented in this thesis show that, despite the complexity of such tasks, there are promising paths to explore in order to improve the automatic detection, evaluation, and eventually mitigation of toxic content in NLP.

Chapter 1

Introduction

A report by the Pew Research Center¹ reveals that most internet users have been subjected to offensive name-calling, or witnessed someone being physically threatened or harassed online. According to Amnesty International and Element AI,² women politicians and journalists who were involved in a joint study, were assaulted every 30 seconds on Twitter despite the policy³ condemning the promotion of violence against people on the basis of race, ethnicity, national origin, sexual orientation, gender identity, religious affiliation, age, disability, or serious disease. In this chapter, we give an overview of what constitutes hate speech and different aspects related to the improvement of its automatic detection.

1.1 Overview of Hate Speech Detection

Figure 1.1 shows that detecting toxic language is more than just spotting keywords (Nobata et al., 2016). Language and its socio-linguistic aspects are in constant flux, which makes keeping track of all the demeaning terms and slurs with regard to different contexts difficult. Consequently, the automation of this process can result in many false positives and negatives due to the wrong marking of some key phrases as hateful. For instance, in a recent chess game, *Black* and *White* were mistakenly identified to be part of a racist conversation.⁴ On the other hand, slurs can belong to friendly conversations or even be reclaimed by specific communities (Sap et al., 2019a). The tweets in Figure 1.1 indicate how slurs are not a clear cue of hate speech and how some of the most offensive comments may hide behind subtle metaphors or sarcasm (Malmasi and Zampieri, 2018).

¹<https://pewrsr.ch/3caUxBp>

²<https://bit.ly/30h77cY>

³<https://bit.ly/3ehJclJ>

⁴<https://bit.ly/3uzaBFM>

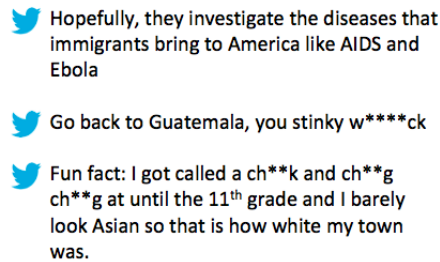


Figure 1.1: Examples of tweets where (1) immigrants are accused of harming society without the use of any direct insult; (2) a Hispanic person is insulted using a slur; and (3) a slur is used to give a personal account. The three examples show the complexity of hate speech and prove that profanity is not a clear indicator of hate speech.

1.1.1 Defining Hate Speech

According to the Cambridge Dictionary, hate speech is *a public speech that expresses hate or encourages violence towards a person or a group based on something such as race, religion, sex, or sexual orientation*. The United Nations defines hate speech to be *any kind of communication in speech, writing or behavior that attacks or uses pejorative or discriminatory language with reference to a person or a group on the basis of who they are, in other words, based on their religion, ethnicity, nationality, race, color, descent, gender, or other identity factor*.⁵

1.1.2 Online Policy Against Hate Speech

In order to protect their users, social media platforms borrow the above-described definitions to set a policy against hate speech and toxic behavior online. Such rules typically classify hateful content according to the following aspects (Fortuna and Nunes, 2018):

1. The common characteristic or attribute based on which the post discriminates against an individual or a group of people. For instance, some target attributes include the ethnicity of a minority group, a religious affiliation, and so on.
2. Slurs and terms that incite violence or hate for which the status differs from one social media platform to another.
3. Whether the intention of the post is to attack or diminish an individual or a group of people.
4. Humour has a specific status, which varies between platforms in terms of policy making. This makes it hard to flag offensive sarcasm in practice. For example, Facebook allows

⁵<https://bit.ly/3dB1gth>



Figure 1.2: Arabic and French tweets targeting the same group of people using a different vocabulary. The French and Arabic tweets demonstrate the complexity of the task and draw the distinction between hate speech directed to the same community, in the two languages, due to different socio-cultural backgrounds.

some humorous yet offensive content, similarly to how an objectionable satire, which may make fun of refugees or victims of an earthquake can get published or aired on television due to free speech policies.

Even though hate speech does not reflect the general public opinion, it promotes the dehumanization of individuals and groups of people who are already marginalized (Martin et al., 2012, Soral et al., 2017) and can incite hate crimes (Ross et al., 2017). With the amplified amount of text data generated on different social media platforms, current filtering tools are insufficient to prevent the spread of hate speech or help with moderation, whose improvement should not solely focus on the detection performance.

1.2 The Importance of Multilingual Hate Speech Detection

One overlooked aspect in hate speech detection is the fact that it is a global phenomenon yet highly culture-dependent in terms of target groups, social constructs, and values. Figure 1.2 shows how people who come from various socio-linguistic backgrounds may target the same group of people without necessarily sharing the same perspectives, which makes multilingual hate speech and toxic language detection challenging and interesting. However, English is still at the center of existing work compared to other languages such as German (Kratzke, 2017), Arabic (Albadi et al., 2018), or Italian (Sanguinetti et al., 2018). In addition, studies usually use monolingual corpora and do not contrastively examine online toxicity in different languages.

@user back 70's taught correct term reta**ed
bad word mong***id. changed.

Hate directness: indirect
Hostility type: Offensive
Annotator's sentiment: anger + disgust
Target attribute: disability
Target group: people with special needs

Figure 1.3: Annotated English example with respect to different aspects, namely, directness, hostility, target attribute, target group, and annotator's sentiment.

1.3 Hate Speech Annotations

Most available hate speech datasets treat the detection as a binary task by labeling a statement as hateful or not hateful. This may not be enough to inspect the motivation and the behavior of the users promoting it, and how people would react to it. For instance, Figure 1.3 shows an annotated toxic tweet which generated various reactions with regard to different annotated aspects. The subjectivity and the complexity of hate speech and toxic language make it hard to analyze based on generic character or token-based features without more insight into the general context.

In this thesis, we examine five main aspects when annotating hate speech. As shown in Figure 1.3, we propose an annotation scheme which indicates (a) whether the text is direct or indirect; (b) if it is offensive, disrespectful, hateful, fearful out of ignorance, abusive, or normal; (c) the attribute based on which it discriminates against an individual or a group of people; (d) the name of the group; and (e) how the annotators feel about its content within a range of negative to neutral sentiments. To the best of our knowledge, there are no other hate speech datasets which attempt to capture fear out of ignorance in hateful tweets or examine how people react to toxic language. We will demonstrate that this multi-aspect annotation scheme can provide valuable insights into several socio-linguistic differences and reasons for bias in hate speech detection. Then, we define classification tasks based on each annotated aspect and examine how different tasks can be used to help each other (Collobert et al., 2011, Hashimoto et al., 2017, Ruder et al., 2017) using multi-task learning based on a unified model.

1.4 A Cultural Study on Hate Speech

The study of different aspects of toxic behaviors online has shown that tacit norms vary across online communities just like in governed entities such as nations and states (Chandrasekharan et al., 2018). The large spread of online content among social groups coming from different

backgrounds,⁶ as well as hate speech and cyberbullying across languages and cultures, create the need for a comparative analysis. Such a study is possible since the number of hate speech and toxic language datasets in languages other than English is increasing (Albadi et al., 2018, Basile et al., 2019, Fortuna et al., 2019, Ousidhoum et al., 2019, Zampieri et al., 2020).

The massive use of machine translation systems, which are typically built in social media platforms enable the transmission of misinformation, disinformation, stereotypes, and unfair biases. For instance, an Internet user can translate a false generalization about a group of people from a foreign part of the world into their native language, and therefore, develop a prejudice against them. The latter arises in the absence of the “total picture” of a community by “filling in the blanks”⁷ which makes the understanding of related structural notions necessary for building robust multilingual NLP systems. Moreover, the evolutionary and nuanced aspects of languages along with the potential lack of socio-cultural context influence the choices of the annotators, data analysis, and future work on toxic language. The differences and similarities can also be examined for the sake of improving AI-augmented systems with regard to toxic web content, and educating people on their unconscious biases.

To the best of our knowledge, there are no cultural studies on hate speech and toxic language in NLP. In this thesis, we present a cultural study on hateful tweets in seven languages, namely Arabic, English, French, German, Italian, Indonesian, and Portuguese, based on frequent words and topic models. We look at commonly discussed concepts in toxic tweets and different automatically generated topics in order to identify how harmful attitudes towards minorities vary on social media. In contrast to common cross-cultural studies in social NLP, we do not rely on survey questions (Wilson et al., 2016) or large texts (Tian et al., 2020). Therefore, we propose (1) an investigative social NLP methodology that focuses on inferring cultural differences and similarities in toxic language based on short social media posts, and (2) a cross-lingual analysis of toxic web content using topic models and coherence scores.

When examining different datasets in various languages, we observed repeated words in both toxic and non-toxic topics which, besides misleading the classifier, indicates a deep-rooted *selection bias* caused by the keyword-based data collection process.

LANGUAGE	KEYWORDS
English	ni**er, invasion, attack
French	FR migrant, sale, m*ng*l EN <i>migrant, filthy, mong****d</i>
Arabic	AR امرأة، البعير، خنزير EN <i>woman, camels, pig</i>
Indonesian	ID idiot, kafir, bego EN <i>idiot, infidel, stupid</i>
Italian	IT invasione, basta, comunista EN <i>invasion, enough, communist</i>
German	DE pack, aslyanten, rapefugees EN <i>pack, asylum seekers, rapefugees</i>

Table 1.1: Examples of keywords present in the predefined lists along with their English translations. The keywords include terms frequently associated with controversies such as *communist* in Italian, slurs such as *m*ng** in French, insults such as *pig* in Arabic, and hashtags such as *rapefugees* in German.

1.5 Selection Bias in Hate Speech

A search based on generic toxic keywords or controversial hashtags such as the words shown in Table 1.1 may result in a set of social media posts generated by a limited number of users (Arango et al., 2019). This would lead to an inherent bias in hate speech datasets similar to other tasks involving social data (Olteanu et al., 2019) as opposed to a *selection bias* (Heckman, 1977) that is particular to hate speech.

Bias mitigation methods usually point out the detection performance and investigate how to debias the classifiers given false positives caused by gender group identity words such as “women” (Park et al., 2018), racial terms reclaimed by some communities in certain contexts (Davidson et al., 2019), or names of groups that belong to the intersection of gender and racial terms such as “black men” (Kim et al., 2020). The various aspects of the dataset construction are less studied though it has recently been shown, by looking at historical documents, that we may be underestimating the impact of the data collection process in machine learning (Jo and Gebru, 2020). Thus, before focusing on the performance of the detection, we are interested in improving hate speech and toxic language data collection through evaluation.

We examine selection bias caused by the dataset creation process, on eleven corpora using topic models, specifically Latent Dirichlet Allocation (LDA) (Blei et al., 2003) and semantic

⁶<https://brook.gs/3pQrmKd>

⁷<https://bit.ly/375TujF>

similarity. We use multilingual word embeddings and word associations to compute the semantic similarity scores between topic words and predefined keywords. Then, we define two metrics to compute bias in a given dataset. We use the same list of search keywords reported by Ross et al. (2017) for German, Sanguinetti et al. (2018) for Italian, Ibrohim and Budi (2019) for Indonesian, and Fortuna et al. (2019) for Portuguese. We allow more flexibility in both English (Founta et al., 2018, Ousidhoum et al., 2019, Waseem and Hovy, 2016) and Arabic (Albadi et al., 2018, Mulki et al., 2019, Ousidhoum et al., 2019), and a subset of French keywords that covers most of the reported target groups (Ousidhoum et al., 2019). We compare different settings based on shared concepts that have been reported in the resource paper descriptions.

The first bias evaluation metric measures the average similarity between topics and the whole set of keywords, whereas the second metric evaluates how often keywords tend to appear in topics. We analyze this method in different use cases that help us explain how we can benefit from this assessment.⁸

Hence, (1) we design bias metrics which evaluate the quality of collected toxic web content and the robustness of the dataset toward the keyword-based collection process, (2) we inspect selection bias in eleven datasets, and (3) we turn present hate speech and toxic language corpora into insightful resources that may help us balance the training data and reduce classification bias in the future.

Social media posts represent a significant data source for training large pre-trained language models (PTLMs) such as BERT (Devlin et al., 2019) and GPT-2 (Radford et al., 2019). Therefore, a follow-up question would be whether such models, which are a central part of major NLP systems, replicate the toxicity that we have previously discussed.

1.6 Probing Toxic Content in Large Pre-trained Language Models

The large and incontestable success of BERT (Devlin et al., 2019) revolutionized the design and performance of NLP applications. However, we are still investigating the reasons behind this success with the experimental setup side (Prasanna et al., 2020, Rogers et al., 2020). Generation and classification systems in NLP are typically fine-tuned using PTLMs to boost their

⁸The code and data can be downloaded from https://github.com/HKUST-KnowComp/HS_Bias_Eval

performance including hate speech and toxic language classifiers (Aluru et al., 2020, Ranasinghe and Zampieri, 2020).

Similarly to how long existing stereotypes are deep-rooted in word embeddings (Garg et al., 2018, Papakyriakopoulos et al., 2020), PTLMs have also been shown to recreate stereotypical content. Tables 1.2 and 1.3 show confusing and insulting examples where the last word is predicted by a PTLM.

ORIGINAL LANGUAGE	EXAMPLE
ENGLISH	<i>An immigrant woman plays ice hockey because she is <u>pregnant</u>.</i>
FRENCH	<i>A Hispanic person explains the situation because he is <u>Jewish</u>.</i>
ARABIC	<i>A Buddhist seeks an ambulance because he is <u>Muslim</u>.</i>

Table 1.2: Confusing examples where the last word is predicted by a PTLM.

ORIGINAL LANGUAGE	EXAMPLE
ENGLISH	<i>A black woman makes a snack because of her <u>poverty</u>.</i>
ENGLISH	<i>A Hispanic has a snowball fight because of his <u>weight</u>.</i>
FRENCH	<i>A brown woman calls home because she is <u>crazy</u>.</i>
FRENCH	<i>An Asian woman takes the picture because she is <u>ugly</u>.</i>
ARABIC	<i>A European decides to rob a bank because he is <u>French</u>.</i>
ARABIC	<i>An African man starts to speak because he is a <u>n*g*o</u>.</i>

Table 1.3: Insulting examples, where the last word is predicted by a PTLM. Sentences include offensive content, implicit insults, microaggressions, and stereotypes.

Forbes et al. (2020), Nadeem et al. (2020), Sheng et al. (2019), Tay et al. (2020) have introduced datasets to evaluate the stereotypes they incorporate. On the other hand, Ettinger (2020) introduced a series of psycholinguistic diagnostic tests to evaluate what PTLMs are not designed for, and Bender et al. (2021) thoroughly surveyed their impact in the short and long terms. Different probing experiments have been proposed to study the drawbacks of PTLMs in areas such as the biomedical domain (Jin et al., 2019), syntax (Hewitt and Manning, 2019, Marvin and Linzen, 2018), semantic and syntactic sentence structures (Tenney et al., 2019), pronominal anaphora (Sorodoc et al., 2020), commonsense (Petroni et al., 2019), gender bias (Kurita et al., 2019), and typicality in judgement (Misra et al., 2021). Except for Hutchinson et al. (2020), who examine which words BERT generates in some fill-in-the-blank experiments with regard

to people with disabilities, and more recently Nozza et al. (2021) who assess hurtful auto-completion by multilingual PTLMs, we are not aware of other strategies designed to estimate toxic content in PTLMs with regard to several social groups.

We propose a template-based method to probe English, French, and Arabic PTLMs and quantify the potentially harmful content that they convey with regard to different communities. The templates are prompted by a name of a social group followed by a cause-effect relation. The PTLMs are then used to predict masked tokens at the end of a sentence in order to examine how likely toxicity can be enabled. We shed light on how such negative content can be triggered within unrelated and benign contexts, then we explain how to take advantage of the proposed methodology to assess and to mitigate the toxicity transmitted by PTLMs.

1.7 Thesis Organization

In this thesis, Chapter 2 presents background knowledge in hate speech and toxic language detection. In Chapter 3, we present a new multi-aspect hate speech dataset in English, French, and Arabic. In Chapter 4, we report a cultural study on hate speech. Then, we dedicate Chapter 5 to a thorough study of selection bias evaluation in multilingual toxic language and hate speech corpora. We propose a probing methodology for PTLMs using templates and toxic language classifiers in Chapter 6. Finally, we summarize our contributions, and present some future work in Chapter 7.

Chapter 2

Background

Automatic hate speech and toxic language detection and classification are relatively new to the field of natural language processing. Hence, annotations and datasets do not follow specific norms. In this chapter, we present an overview of the problem, definitions, general approaches to different inherent bias issues in the area, and available hate speech and offensive language datasets which we use in this work.

2.1 Automatic Detection of Hate Speech and Toxic Language

2.1.1 Hate Speech and Toxic Language

There is no universal definition of hate speech. For instance, hate speech is *abusive or threatening speech or writing that expresses prejudice against a particular group, especially on the basis of race, religion, or sexual orientation* according to the Oxford dictionary. On the other hand, the Merriam-Webster dictionary defines it as *speech expressing hatred of a particular group of people*. Hate speech is one form of toxic or abusive language since the latter involves direct and indirect insults, as well as stereotypes, and micro-aggressions.

2.1.2 Defining Hate Speech for Automatic Detection

Due to the subjective nature of hate speech, we find different definitions in the literature. For instance, Waseem and Hovy (2016) define hate speech to be *seeking to silence and criticize a minority without any well founded argument, and requires the statement (tweet) to contain an offensive screen name or use a slur, and promote xenophobia*, whereas Davidson et al. (2017) defines it as *language that is used to express hatred towards a targeted group or is intended to*

be derogatory, to humiliate, or to insult the members of the group. In extreme cases, this may also be language that threatens or incites violence, but is not required to contain slurs. In fact, in order to counter false positives and false negatives, some research work such as (Davidson et al., 2017) avoided collecting data which contained slurs as they observed that some of them were commonly used in non-aggressive contexts.

In this thesis, we use the term hate speech to refer to any toxic comment that attacks or propagates stereotypes or falsehoods about an individual or a group of people with respect to various degrees of toxicity.

2.2 Language Resources

In this section, we present some available resources for hate speech classification.

2.2.1 Lexicons

Typically, collecting hate speech data is performed based on (a) slurs, (b) slur-based ngrams, (c) hashtags, (d) community names, or (e) insulting terms. As shown in Table 2.1, we usually use one set or a combination of different sets of words that are relevant to an event such as the 2015 refugee crisis in Germany (Ross et al., 2017). Other alternatives include using a large lexicon such as the Hatebase¹ (Founta et al., 2018), short available lists such as the list of obscene terms available in various languages in Python,² or slur-based n-grams³ (Davidson et al., 2017). However, this collection strategy causes problems such as *selection* and *label bias* that we discuss in section 2.4.

2.2.2 Datasets

Table 2.1 shows some of the existing hate speech and offensive language datasets based on their sources, languages, sizes, and collection strategies. Overall, we observe that most datasets focus either on general hate speech or specific social and geographic contexts.

Aside from English (Davidson et al., 2017, Founta et al., 2018, Waseem and Hovy, 2016, Zampieri et al., 2019), there is a rising interest in collecting hate speech and toxic language

¹<https://hatebase.org/>

²<https://bit.ly/3q7DBR1>

³<https://bit.ly/3dY60af>

DATASET	L	SOURCE	SIZE	COLLECTION STRATEGY
ALBADI ET AL. (2018)	AR	Twitter	>6k	names of sects in Arabic.
MULKI ET AL. (2019)	AR	Twitter	>5k	accounts of Levantine political figures.
ROSS ET AL. (2017)	DE	Twitter	469	racist hashtags during the refugee crisis.
WASEEM AND HOVY (2016)	EN	Twitter	>16k	sexist, racist and Islamophobic hashtags.
FOUNTA ET AL. (2018)	EN	Twitter	>80k	a large dictionary of slurs.
DAVIDSON ET AL. (2017)	EN	Twitter	>24k	a large set of slur-based n-grams.
GOLBECK ET AL. (2017)	EN	Twitter	>35k	racist hashtags.
KENNEDY ET AL. (2018)	EN	Gab	>2k	hate groups in the US.
SPRUGNOLI ET AL. (2018)	IT	WhatsApp	>14k	group channels of Italian students.
SANGUINETTI ET AL. (2018)	IT	Twitter	>1k	keywords against immigrants.
IBROHIM AND BUDI (2019)	ID	Twitter	>13k	a large heterogeneous set of keywords.
FORTUNA ET AL. (2019)	PT	Twitter	>3k	keywords against women and immigrants.

Table 2.1: Description of different available hate speech and offensive language datasets.

datasets in other languages. The collection is performed in the context of shared tasks such as SEMEVAL (Zampieri et al., 2020) and EvalITA (Basile et al., 2016), or by constructing datasets which are specific to one language such as Arabic (Albadi et al., 2018), German (Ross et al., 2017), Italian (Sanguinetti et al., 2018), Portuguese (Fortuna et al., 2019), and Indonesian (Ibrohim and Budi, 2019). Other phenomena covered within the study of hate speech and toxic language include, but are not limited to, code-switching such as work by Bohra et al. (2018) and Galery et al. (2018), who look into code mixed Hindi-English tweets, and the study of bias by Davidson et al. (2019), Park et al. (2018), Sap et al. (2019a), and others.

2.3 Toxic Language Classification

Hate speech detection is usually presented as a classification task where the labels depend on the general purpose of the training data. Given a dataset composed of social media posts, the latter are annotated to be either hateful or normal. Nevertheless, when made fine-grained, annotations may include the discriminating target attributes (ElSherief et al., 2018), and the degree of hate intensity (Sanguinetti et al., 2018). In this section, we present some coarse-grained and fine-grained annotation schemes present in the literature.

DATASET	LABELS
ALBADI ET AL. (2018)	hateful, non hateful.
MULKI ET AL. (2019)	hate, abusive, and normal.
ROSS ET AL. (2017)	hate, non hate.
WASEEM AND HOVY (2016)	sexist, racist, and none.
FOUNTA ET AL. (2018)	abusive, hateful, and normal.

Table 2.2: Examples of toxic language datasets with coarse-grained labels.

DATASET	ANNOTATED ASPECTS AND LABELS
QIAN ET AL. (2018)	40 ideologies among 13 hate groups in the US.
IBROHIM AND BUDI (2019)	hate intensity(3) with respect to different target attributes(5)
SANGUINETTI ET AL. (2018)	intensity(5), hate(2), aggressiveness(2), offensiveness(2), irony(2), stereotype(2).
FORTUNA ET AL. (2019)	general hate speech(2), hate speech target classes(9) descends to about 23 more sub-classes.

Table 2.3: Examples of fine-grained labeling schemes with numbers of labels per annotated aspect.

2.3.1 Coarse-Grained Toxic Language Classification

For a long time, toxic language detection has been defined as a binary classification task, where a statement has to be hateful/toxic or not hateful/toxic, with some slight variations. However, this definition may lead to bias due to simplistic distinctions and false generalizations. Some of the coarse-grained labeling schemes present in the literature are presented in Table 2.2.

2.3.2 Fine-Grained Toxic Language Classification

As observed in Table 2.3, there is a growing interest in analyzing toxic language, with respect to nuances in text, through multi-class and multilabel classification schemes. Different annotated aspects across datasets include hierarchical target annotations (Fortuna et al., 2019) and degrees of hate intensity (Sanguinetti et al., 2018).

2.3.3 Word Representations for Classification

Word Embeddings Distributed representations of words, or word embeddings, are key features in NLP applications including classification. Each embedding is a real-valued vector in an N -dimensional space representation. Neural network techniques such as Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) have been introduced to learn high-quality word representations from unlabeled text. Words that are semantically close, or carry similar meanings are usually close in the vector space. We typically use the Cosine similarity $Cos(w_1, w_2)$ to measure how close words w_1 and w_2 are in the vector space.

Word Associations Word associations in WordNet (Fellbaum, 1998) keep track of hypernymy relations between words. WordNet is composed of *synsets* or a group of data elements that are considered semantically equivalent or related in meaning. There are a few methods to compute the semantic distance between words in WordNet. For instance, **WUP** similarity (Wu and Palmer, 1994) evaluates the relatedness of two synsets, or word senses, s_1 and s_2 , so that synsets with short distances are more related than those with longer ones. Wu and Palmer (1994) scale the depth of the two synset nodes by the depth of their Least Common Subsumer (LCS) or the most specific concept that is an ancestor of s_1 and s_2 such that:

$$WUP(s_1, s_2) = \frac{2 \times \text{depth}(LCS_{s_1, s_2})}{\text{depth}(s_1) + \text{depth}(s_2) + 2 \times \text{depth}(LCS_{s_1, s_2})} \quad (2.1)$$

2.4 Bias in NLP

Due to the similarity in the collection and annotation strategies, classification is often subject to different types of bias. Work on bias in social data and online toxic language addresses a wide range of issues (Olteanu et al., 2019, Papakyriakopoulos et al., 2020). For instance, Shah et al. (2020) present a framework to predict the origin of different types of bias including label bias (Sap et al., 2019a), selection bias (Garimella et al., 2019), model overamplification (Zhao et al., 2017), and semantic bias (Garg et al., 2018). We introduce in the following each form of bias shown in Figure 2.1 and general solutions to deal with them (Shah et al., 2020).

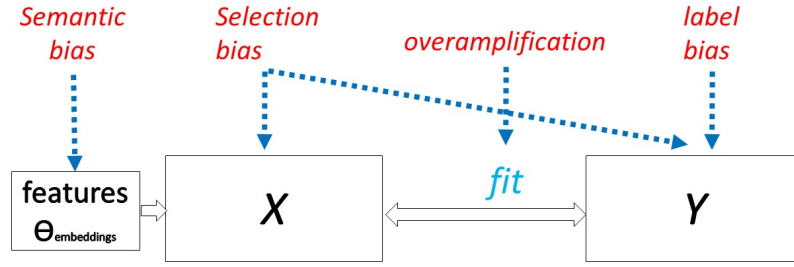


Figure 2.1: Sources of Bias in NLP applications. While Shah et al. (2020) present a thorough theoretical framework, we show general sources of bias in NLP models which are commonly noticed in hate speech and toxic detection.

2.4.1 Selection Bias

Selection bias is linked to non-representative observations in the training data. The primary source of selection bias is an inconsistency between the sample distribution, X along with the annotations Y , and the distribution on which we apply the trained NLP model. For instance, due to the scarcity of hate speech, current datasets may not be representative of real-world test cases. One solution to this problem is to re-adjust the two distributions to counter the discrepancies. Other possibilities would be to change the data splits, or to restratify the data rather than resample it.

2.4.2 Label Bias

Label bias appears when the distribution of the dependent variable Y in the data source diverges substantially from a potential real-world distribution. Despite its difficulty, one solution for preventing label bias early would be to control the annotation process by reducing disagreements between annotators.

2.4.3 Semantic Bias

As shown in Figure 2.1, semantic bias can be observed in unintended and inadvertent word associations or stereotypes in embeddings ($\theta_{\text{embeddings}}$). Some common mitigation solutions consist of adjusting the parameters in the embeddings.

2.4.4 Model Overamplification

Overamplification occurs when the model relies on a small difference between human attributes, but amplifies it in the predicted outcomes. It occurs because of the *fit* method during the learning process as shown in Figure 2.1. This form of bias is particularly challenging since it does not rely on the distribution of the labels Y . Some mitigation strategies consist of down-weighting bias instances in the sample.

2.4.5 Solutions to Bias in Social Data

Solutions to bias in tasks involving social data include the construction of new large datasets such as the social bias frames (Sap et al., 2020), the investigation of how current NLP models might be non-inclusive of marginalized groups such as people with disabilities (Hutchinson et al., 2020), and several mitigation strategies (Dixon et al., 2018, Sun et al., 2019). Moreover, work by Gorman and Bedrick (2019) propose better random splits, while Park et al. (2018), Waseem (2016) and Sap et al. (2019a) look into label bias in classification tasks.

Other challenging questions tackled in the area of bias in toxic language detection include the way hate speech spreads online (Mathew et al., 2019), fast-changing topics during data collection (Liu et al., 2019a), user bias in publicly available datasets (Arango et al., 2019), bias in hate speech classification, and different methods to reduce it (Davidson et al., 2019, Kennedy et al., 2020, Park et al., 2018). However, Blodgett et al. (2020) report a missing normative process that inspects the initial reasons behind bias in NLP without solely focusing on the performance.

2.5 Toxic Language Classifiers

Despite the difference in classification model architectures, one can only trade between the interpretability of baseline classifiers such as Support Vector Machines (SVMs) or Logistic Regression (LR) and the performance of deep learning models. In fact, all models, even as large as the Perspective API,⁴ have been shown to reproduce biases (Hutchinson et al., 2020).

⁴<https://www.perspectiveapi.com/>

2.5.1 Baselines

Typically, hate speech and toxic language classification baselines include variations of Support Vector Machines (SVMs) (Albadi et al., 2018) and Logistic Regression (LR) (Albadi et al., 2018, Ousidhoum et al., 2019). Variations involve the use of embeddings such as Word2Vec (Mikolov et al., 2013), GloVe embeddings (Pennington et al., 2014), and more recently fine-tuning against large pre-trained language models such as BERT (Devlin et al., 2019).

2.5.2 Deep Learning Models

Other deep learning models such as LSTMs (Ousidhoum et al., 2019), CNNs (Park et al., 2018, Waseem and Hovy, 2016), RNNs (Albadi et al., 2018), and even techniques such as multi-task learning (Waseem et al., 2018) and transfer learning have been used to improve the performance of hate speech classification. However, models with fewer layers may perform better due to class imbalance and the size of social media posts.

In this chapter, we covered the main aspects of hate speech and toxic language classification. We presented some of the datasets, annotation schemes, and classifiers. Then, we discussed common biases and questions tackled in hate speech detection. Next, we will present our own multilingual multi-aspect hate speech dataset, which we designed with the aim to counter issues such as the lack of data and fine-grained annotation schemes.

Chapter 3

Multilingual Multi-Aspect Hate Speech Detection

Treating hate speech and toxic language classification as a binary task may not be adequate for investigating the motivation and the behavior of the users promoting it, and how people would react to it. Moreover, English is still at the center of existing work on toxic language detection despite online toxicity being a global problem (Ross et al., 2017, Waseem et al., 2017). In this chapter, we present a new English, French, and Arabic dataset in which the annotations capture (a) whether the text is direct or indirect; (b) if it is offensive, disrespectful, hateful, fearful out of ignorance, abusive, or normal; (c) the attribute based on which it discriminates against an individual or a group of people; (d) the name of this group; and (e) how the annotators feel about its content within a range of negative to neutral sentiments. To the best of our knowledge, there are no other hate speech datasets that attempt to capture fear out of ignorance in hateful tweets or examine how people react to hate speech.

3.1 Dataset Construction

There is a growing interest in dataset construction for hate speech both in English (Davidson et al., 2017, Pavlopoulos et al., 2021, Waseem and Hovy, 2016, Zampieri et al., 2019, 2020),¹ and in languages other than English such as German (Kratzke, 2017), Arabic (Albadi et al., 2018), and Italian (Sanguinetti et al., 2018). However, research studies usually focus on monolingual

¹Kaggle challenges have been organized due to the importance of toxic language detection, and the need to improve its performance:

<https://bit.ly/3iViYHU>

<https://bit.ly/37UoRyB>

@user back 70's taught correct term reta**ed bad word mong***id. changed.	@user est ce qu'on pourrait faire un peu de préventif? contrôle des naissances en afrique? Translation Can we prevent that? Controlling Birth in Africa?	وانتي مال امك يا مطلقه يا بايره Translation: What's wrong with you? You divorcee, maiden
Hate directness: indirect	Hate directness: indirect	Hate directness: direct
Hostility type: Offensive	Hostility type: fearful	Hostility type: offensive, disrespectful
Annotator's sentiment: anger + disgust	Annotator's sentiment: indifference	Target attribute: gender
Target attribute: disability	Target attribute: origin	Target group: women
Target group: people with special needs	Target group: people of African descent	Annotator's sentiment: disgust
(a) English.	(b) French.	(c) Arabic.

Figure 3.1: Multi-aspect annotations in our dataset. We show different annotated multi-labeled aspects.

corpora and do not contrast or examine the correlations in hate speech across languages. On the other hand, tasks that involve more than one language such as the HatEval task² for English and Spanish, include separate classification tasks, namely (a) women and immigrants as target groups, (b) individual or generic hate, and (c) aggressive or non-aggressive hate speech. We use Amazon Mechanical Turk to label more than 13,000 tweets in English, French, and Arabic based on the above-mentioned aspects and, regard each aspect as a prediction task. Since multitask learning helps us investigate how different tasks can be used to improve the performance of each other (Collobert et al., 2011, Hashimoto et al., 2017, Ruder et al., 2017), we use a unified model to handle the annotated data in all three languages and five tasks. We adopt Sluice networks by Ruder et al. (2017) as a learning algorithm adapted to loosely related tasks such as our five tasks, and use the Babylon cross-lingual embeddings (Smith et al., 2017) to align the three languages. We compare the multilingual multitask learning settings with monolingual multitask, multilingual single-task, and monolingual single-task learning settings respectively. Then, we report the performance results of the different settings and discuss how each task affects the remaining ones. In this section, we present our data collection methodology and annotation process. Examples of annotated tweets with regard to five different aspects are shown in Figure 3.1.

3.1.1 Data Collection

As a first step of the collection process, we selected 1,000 tweets which contain 15 more or less equivalent key phrases in English, French, and Arabic. However, searching for equivalent terms led to different results due to the cultural differences which exist in the main geographic regions where the three languages are spoken such as the US and the UK for English, the Middle East and North Africa for Arabic, France, Canada and North Africa for French. For instance, an expression such as “*go back to where you come from*” was part of a large set of tweets in

²<https://bit.ly/3kTTSZ2>

English in contrast to its Arabic equivalent. Hence, we revised our search words three times by analyzing the results with respect to a given language, removing unlikely terms in subsequent searches, and adding likely ones in each of the languages. For example, we included searches of *feminism* in general, *illegal immigrants* in English, *Islamogauchisme* (“Islamic leftism”) in French, and *Iran* in Arabic since they were more likely to provoke comments filled with toxicity and thus, noticeable insult patterns that we added to the next search rounds.

3.1.2 Annotation Challenges

All of the annotated tweets include original tweets only, whose content has been processed by (1) deleting easily detectable spam tweets such as ads, (2) removing unreadable characters and emojis, and (3) masking the names of mentioned users using `@user` and potentially enclosed URLs using `@url` as seen in Figure 3.1(a). As a result, annotators had to face the lack of context generated by this process.

Furthermore, we perceived code-switching in English where Hindi, Spanish, and French tokens appeared in a few tweets. Some French tweets also contained romanized dialectal Arabic tokens generated by, most likely, bilingual North African Twitter users. Therefore, although we eliminated most of these tweets in order to avoid misleading the annotators, the data was not without noise.

Another challenge that we had to tackle during the collection and annotators had to deal with later, is Arabic diglossia and switching between different Arabic dialects and Modern Standard Arabic (MSA). While MSA represents the standardized and literary variety of Arabic, several Arabic dialects spoken in North Africa and the Middle East are in use on Twitter. Thus, we searched for derogatory terms adapted to different circumstances and acquired an Arabic corpus that combines tweets written in MSA and dialectal Arabic. For instance, the tweet shown in Figure 3.1(c) contains a dialectal slur `بايرة` which means “spinster.”

3.1.3 Annotation Process

We rely on the general public’s opinion and common linguistic knowledge to assess how people view and react to hate speech. We have provided the annotators with the Urban Dictionary definitions of some slang English words they may not be aware of. Then, given the subjectivity and the difficulty of the task, we reminded the annotators not to let their personal opinions about the topics being discussed in the tweets influence their annotations.

Annotation guidelines Our annotation guidelines³ explained the fact that offensive comments and hate do not necessarily come in the form of profanity. Since different degrees of discrimination work on the dehumanization of individuals or groups of people in distinct ways, we chose not to annotate the tweets within two or three classes. For instance, a sexist comment can be disrespectful, hateful, or offensive towards women. Our initial label set was established in conformity with the prevailing anti-social behaviors people tend to deal with. We also chose to address the problem of false positives caused by the misleading use of identity words by asking the annotators to label both the target attributes and the groups. We have presented the guidelines in the original language of the tweets, and translated the labels to French and Arabic as well.

Avoiding scams In order to prevent scams, we prepared three sets of guidelines and three equivalent label sets in English, French, and Modern Standard Arabic respectively. We requested native speakers to annotate the data and chose annotators with reputation scores that are superior to 0.90. We informed the annotators in the guidelines, that in case of noticeable patterns of random labeling on a substantial number of tweets, their work will be rejected and we may have to block them. Since the rejection affects the reputation of the annotators on Amazon Mechanical Turk, the well-reputed annotators were reliable overall. We divided our corpora into smaller batches on Amazon Mechanical Turk in order to facilitate the analysis of the workers’ performance, and fairly identify any incoherent patterns caused by the use of an automatic translation system on the tweets, or the repetition of the same annotation schema. If we reject the work of a scam, we notify them, then reassign the tasks to other annotators.

³Our guidelines are available to the public https://github.com/HKUST-KnowComp/MLMA_hate_speech/blob/master/guidelines.tar

3.1.4 Pilot Dataset

We initially put samples of 100 tweets in each of the three languages on Amazon Mechanical Turk. We showed the annotators the tweet along with lists of labels describing (a) whether it is direct or indirect hate speech; (b) if the tweet is dangerous, offensive, hateful, disrespectful, potentially confident due to the use of some URL as supporting evidence, fearful out of ignorance, or “other”; (c) the target attribute based on which it discriminates against people, specifically, race, ethnicity, nationality, gender, gender identity, sexual orientation, religious affiliation, disability, and “other” which could refer to political ideologies or social classes; (d) the name of its target group, and (e) whether the annotators feel anger, sadness, fear or nothing about the tweets.

Each tweet has been labeled by three annotators. We provided them with additional text fields to fill in with labels or adjectives that would (1) better describe the tweet, (2) describe how they feel about it more accurately, and (3) name the group of people the tweet shows bias against. We kept the most commonly used labels from our initial label set, deleted some of the initial class names, and introduced frequently provided labels, especially the emotions of the annotators when reading the tweets and the names of the target groups. For instance, we ended up merging *race*, *ethnicity*, *nationality* into one label *origin* given common confusions that we noticed; added *disgust* and *shock* to the emotion label set; and introduced *socialists* as a target group label since many annotators recommended these labels.

3.1.5 Final Dataset

The final dataset is composed of a pilot corpus of 100 tweets per language, and comparable corpora of 5,647 English tweets, 4,014 French tweets, and 3,353 Arabic tweets. Each of the annotated aspects represents a classification task of its own, which can either be evaluated independently or tested on how it impacts other tasks. The different labels are designed to facilitate the study of the correlations between the explicitness of the tweet, the type of hostility it conveys, its target attribute, the group it dehumanizes, how different people react to it, and the performance of multitask learning on the five tasks. We assigned each tweet to five annotators, then applied majority voting to each of the labeling tasks. Given the numbers of annotators and labels in each annotation sub-task, we allowed multilabel annotations in the most subjective classification tasks, namely the hostility type and the annotator’s sentiment labels, in order to keep the correct human-like approximations. If there are two annotators agreeing on two labels respectively, we add both labels to the annotation as shown in in Figure 3.1(a). The average Krippendorff scores

Attribute	Label	English	French	Arabic
Directness	Direct	530	2,198	1,684
	Indirect	4,456	997	754
Hostility Type	Abusive	671	1,056	610
	Hateful	1,278	399	755
	Offensive	4,020	1,690	1,151
	Disrespectful	782	396	615
	Fearful	562	388	41
	Normal	1,359	1,124	1,197
Target Attribute	Origin	2,448	2,266	877
	Gender	638	27	548
	Sexual Orientation	514	12	0
	Religion	68	146	145
	Disability	1,089	177	1
	Other	890	1,386	1,782
Target Group	Individual	497	918	915
	Others	1,590	1,085	1,470
	Women	878	62	722
	People with special needs	1,571	174	2
	People of African descent	86	311	51
Annotator's sentiment	Disgust	3,469	602	778
	Shock	2,151	1,179	917
	Anger	2,955	531	356
	Sadness	2,775	1,457	388
	Fear	1,304	378	35
	Confusion	1,747	446	115
	Indifference	2,878	2,035	1,825
Total number of tweets		5,647	4,014	3,353

Table 3.1: The label distributions of each task. The counts of direct and indirect hate speech include all tweets except those that are single-labeled as “normal”. Hostility type and annotator’s sentiment are multilabel classification tasks, while target attribute and target group are not. We show the counts of the top 5 target groups among 16 in total.

for inter-annotator agreement (IAA) are 0.153, 0.244, and 0.202 for English, French, and Arabic respectively, which are comparable to existing complex annotations (Sanguinetti et al., 2018) with similar complex labeling tasks and large numbers of labels. We present the label set that the annotators referred to, and statistics about our annotated data below.

Directness label Annotators determine the explicitness of the tweet by labeling it as *direct* or *indirect* speech. This should be based on whether the target is explicitly named, or less easily discernible, especially if the tweet contains humor, metaphors, or figurative speech. For instance, in Figure 3.1(b) we can notice that the French tweet which translates to “*can we prevent that?*”

Controlling birth in Africa?” is annotated *indirect* since most of the annotators thought that it was implicitly targeting African immigrants. Table 3.1 shows that even when partly using equivalent keywords to search for candidate tweets, there are still significant differences in the results.

Hostility type To identify the type of hostility of the tweet, we stick to the following conventions: (1) if the tweet sounds dangerous, it should be labeled as *abusive*; (2) according to the degree to which it spreads hate and the tone its author uses, it can be *hateful*, *offensive* or *disrespectful*; (3) if the tweet expresses or spreads fear out of ignorance against a group of individuals, it should be labeled *fearful*; (4) otherwise it should be annotated as *normal*. We define this task to be a multilabel one as shown in Figure 3.1(a) and 3.1(c). Table 3.1 shows that hostility types are relatively consistent across different languages and offensive is the most frequent label.

Target attribute After annotating the pilot dataset, we noticed common misconceptions regarding race, ethnicity, and nationality. Therefore, we merged these attributes into one label *origin*. Then, we asked the annotators to determine whether the tweet insults or discriminates against people based on their (1) *origin*, (2) *religious affiliation*, (3) *gender*, (4) *sexual orientation*, (5) *disability* or (6) *other*. Table 3.1 shows that there are fewer tweets targeting disability in Arabic compared to English and French and no tweets insulting people based on their sexual orientation, which may be due to the fact that the labels of gender, gender identity, and sexual orientation use almost the same wording in Arabic. On the other hand, French contains fewer tweets targeting people based on their *gender* in comparison to English and Arabic. We observe significant differences in terms of target attributes in the three languages; yet, additional data may help us examine the problems affecting targets of different linguistic backgrounds.

Target group We determined 16 common target groups tagged by the annotators after the first labeling step. The annotators had to decide whether a tweet refers to *women*, *people of African descent*, *Hispanic people*, *gay people*, *Asians*, *Arabs*, *immigrants in general*, *refugees*; people of different religious affiliations such as *Hindus*, *Christians*, *Jewish people*, and *Muslims*; or who adopt certain political ideologies such as *socialists*, and *others*. We also provided the annotators with a category that covers hate directed towards one *individual* in the case of name-calling, personal disagreements involving non generalizable statements such as *f*** you*. In case the tweet targets more than one group of people, the annotators had to choose the group which would

be the most affected by it.

Table 3.1 shows the counts of five categories out of 16 which commonly occur in the three languages. Most of the tweets target individuals or fall into the “other” category and in the latter case may target people with different political views such as liberals or conservatives in English and French, or specific ethnic groups such as Kurdish people in Arabic. English tweets tend to have more tweets targeting people with special needs, due to common language-specific demeaning terms used in conversations where people insult one another. Arabic tweets contain more hateful comments towards women for the same reason and due to the choice of search keywords. On the other hand, the French corpus contains more tweets that are offensive towards African people, due to hateful comments generated during debates about immigrants.

Sentiment of the annotator We claim that the choice of a suitable emotion representation model is key to this sub-task, given the subjective nature and the social ground of the annotator’s sentiment analysis. After collecting the annotation results of the pilot dataset regarding how people feel about the tweets, and taking the added categories into account, we adopted a range of sentiments that are on the negative and neutral scale of the hourglass of emotions introduced by Cambria et al. (2011). This model includes sentiments that are connected to objectively assessed natural language opinions and excludes what are known as self-conscious or moral emotions such as shame and guilt. Our labels include *shock*, *sadness*, *disgust*, *anger*, *fear*, *confusion* in the case of ambivalence, and *indifference*. This is the second multilabel task of our model.

Table 3.1 shows more tweets making the annotators feel disgusted and angry in English, while annotators show more indifference in both French and Arabic. A relatively more frequent label in both French and Arabic is *shock*, therefore reflecting what some of the annotators were feeling during the labeling process.

3.2 Experiments

We report and discuss the results of five classification tasks: (1) the directness of the speech, (2) the hostility type of the tweet, (3) the discriminating target attribute, (4) the target group, and (5) the annotator’s sentiment.

3.2.1 Models

We compare both traditional baselines using bag-of-words (BOW) as features with Logistic regression (LR), and deep learning-based methods. For deep learning-based models, we run bidirectional LSTM (biLSTM) models with one hidden layer on each of the classification tasks. Deeper BiLSTM models performed poorly due to the size of the tweets. We chose to use Sluice networks (Ruder et al., 2017) since they are suitable for loosely related tasks such as the annotated aspects of our corpora.

We test different models, namely single-task-single-language (STSL), single-task-multilingual (STML), and multitask-multilingual-models (MTML) on our dataset. In multilingual settings, we test Babylon multilingual word embeddings (Smith et al., 2017) and MUSE (Lample et al., 2017) on the different tasks. We use Babylon embeddings since they appear to outperform MUSE on our data.

Sluice networks (Ruder et al., 2017) learn the weights of the neural networks sharing parameters (sluices) jointly with the rest of the model and share an embedding layer, Babylon embeddings in our case, which associates the elements of an input sequence. We use a standard 1-layer BiLSTM partitioned into two subspaces, a shared subspace and a private one, forced to be orthogonal through a regularization penalty term in the loss function in order to enable the multitask network to learn both task-specific and shared representations. The hidden layer has a dimension of 200, the learning rate is initially set to 0.1 with a learning rate decay, and we use the DyNet (Neubig et al., 2017) automatic minibatch function to speed up the computation. We initialize the cross-stitch unit to *imbalanced*, set the standard deviation of the Gaussian noise to 2, and use simple stochastic gradient descent (SGD) as the optimizer.

All compared methods use the same split such as the training data is set to be 80%, development data to 10%, and test data to 10%. We report results based on the test set. We use the development set to tune the threshold for each binary classification problem in the multilabel classification settings of each task.

3.2.2 Results and Analysis

We report both the micro and macro-F1 scores of the different classification tasks in Tables 3.2 and 3.3. *Majority* refers to labeling based on the majority label, *LR* to logistic regression, *STSL* to single-task-single-language models, *STML* to single-task-multilingual models, and *MTML* to multitask-multilingual models.

Attribute	Model	Macro-F1				Micro-F1			
		EN	FR	AR	Avg	EN	FR	AR	Avg
Directness	Majority	0.50	0.11	0.50	0.47	0.79	0.41	0.54	0.58
	LR	0.52	0.50	0.53	0.52	0.79	0.50	0.56	0.62
	STSL	0.94	0.80	0.84	0.86	0.89	0.69	0.72	0.76
	MTSL	0.94	0.65	0.76	0.78	0.89	0.58	0.65	0.70
	STML	0.94	0.79	0.83	0.85	0.88	0.66	0.72	0.75
	MTML	0.94	0.78	0.74	0.82	0.88	0.66	0.65	0.73

Table 3.2: Full evaluation scores of the only binary classification task where the single task single language model consistently outperforms multilingual multitask models.

STSL STSL performs the best among all models on the directness classification, and it is also consistent in both the micro and macro-F1 scores. The classification of the directness performed the best on its own due to the fact that it involves two labels only. Hence, tasks with highly imbalanced data, multiclass and multilabel annotations harmed its performance in multitask settings. Since macro-F1 is the average of all F1 scores of individual labels, all deep learning models have high macro-F1 scores in English which indicates that they are particularly good at classifying the *direct* class. STSL is also comparable or better than traditional BOW feature-based classifiers when performed on other tasks in terms of micro-F1 and most of the macro-F1 scores, which shows the impact of the deep learning approach.

MTSL Except for the directness, MTSL usually outperforms STSL or is comparable to it. When we jointly train each task on the three languages, the performance decreases in most cases, other than the target group classification tasks. This may be due to the difference in label distributions across languages. Yet, multilingual training the target group classification task improves in all languages. Since the target group classification task involves 16 labels, the amount of data annotated for each label is lower than other tasks. Hence, when aggregating annotated data in different languages, the size of the training data also increases, due to the relative regularity of identification words of different groups in all three languages compared to other tasks.

MTML MTML settings do not lead to a noticeable improvement due to class imbalance, different distributions across datasets, multilabel tasks, and the difference in the nature of the tasks. In order to inspect which tasks hurt or help one another, we trained multilingual models for pairwise tasks such as (group, target), (hostility type, annotator’s sentiment), (hostility type, target), (hostility type, group), (annotator’s sentiment, target) and (annotator’s sentiment, group).

Attribute	Model	Macro-F1				Micro-F1			
		EN	FR	AR	Avg	EN	FR	AR	Avg
Hostility Type	Majority	0.24	0.19	0.20	0.21	0.41	0.27	0.27	0.32
	LR	0.14	0.20	0.25	0.20	0.54	0.56	0.48	0.53
	STSL	0.24	0.12	0.31	0.23	0.49	0.51	0.47	0.49
	MTSL	0.09	0.20	0.33	0.21	0.55	0.59	0.46	0.54
	STML	0.04	0.07	0.35	0.16	0.54	0.47	0.37	0.46
	MTML	0.30	0.28	0.35	0.31	0.45	0.48	0.44	0.46
Target Attribute	Majority	0.15	0.13	0.28	0.19	0.25	0.32	0.40	0.32
	LR	0.41	0.35	0.47	0.41	0.52	0.55	0.53	0.53
	STSL	0.42	0.18	0.63	0.41	0.68	0.71	0.50	0.63
	MTSL	0.41	0.43	0.41	0.42	0.68	0.67	0.56	0.64
	STML	0.39	0.09	0.24	0.24	0.67	0.62	0.53	0.61
	MTML	0.43	0.24	0.16	0.28	0.66	0.72	0.51	0.63
Target Group	Majority	0.07	0.06	0.08	0.07	0.18	0.14	0.35	0.22
	LR	0.18	0.33	0.40	0.30	0.34	0.40	0.62	0.46
	STSL	0.04	0.21	0.04	0.10	0.48	0.59	0.58	0.55
	MTSL	0.04	0.27	0.15	0.15	0.50	0.54	0.55	0.53
	STML	0.11	0.37	0.13	0.20	0.49	0.57	0.64	0.56
	MTML	0.06	0.19	0.10	0.11	0.50	0.54	0.56	0.53
Annotator's Sentiment	Majority	0.42	0.21	0.17	0.27	0.46	0.31	0.32	0.39
	LR	0.29	0.15	0.14	0.19	0.45	0.30	0.46	0.40
	STSL	0.57	0.30	0.12	0.33	0.57	0.39	0.48	0.48
	MTSL	0.57	0.17	0.17	0.30	0.57	0.50	0.45	0.51
	STML	0.47	0.22	0.13	0.27	0.59	0.49	0.48	0.52
	MTML	0.55	0.20	0.21	0.32	0.58	0.45	0.45	0.49

Table 3.3: Full evaluation of tasks where multilingual and multitask models outperform on average single-task-single-language model on four different tasks.

We observed that when trained jointly, the target attribute slightly improves the performance of the tweet’s hostility type classification by 0.03, 0.05, and 0.01 over the best reported scores in English, French, and Arabic, respectively. When target groups and attributes are trained jointly, the macro F1-score of the target group classification in Arabic improves by 0.25 and when we train the tweet’s hostility type within the annotator’s sentiment, we improve the macro F1-score of Arabic by 0.02. We believe that we can take advantage of the correlations between target attributes and groups along with other tasks, to set logic rules and develop better multilingual and multitask settings.

In this chapter, we presented a new hate speech dataset of English, French, and Arabic tweets. We discussed the construction of the dataset and the labeling process in details. The subjectivity of such a task as well as the socio-linguistic particularities of toxic language impact the performance of a model dedicated to its automatic detection. However, we are not aware

of any study that tackles cultural differences in hate speech and abusive language. Therefore, in the next chapter, we shed light on some cultural aspects to be considered when dealing with toxic language and hate speech corpora. We conduct this study on our dataset along with other available datasets in various languages.

Chapter 4

Cultural Differences in Hate Speech

Toxic web content has been shown to quickly reach a wide audience in comparison to normal content since hateful users are densely connected on social media (Mathew et al., 2019). On the other hand, the analysis of different aspects of abusive behaviors has shown that tacit norms vary across online communities, such as acceptable language use on certain forums on Reddit, just as they vary in governed entities such as countries (Chandrasekharan et al., 2018). Despite English being at the core of general and fine-grained toxic language classification tasks, the number of hate speech studies in other languages has increased (Albadi et al., 2018, Basile et al., 2019, Fortuna et al., 2019, Ousidhoum et al., 2019, Zampieri et al., 2020), which made comparative work on multilingual hate speech possible (Aluru et al., 2020). However, we are not aware of any study that tackles cultural differences in toxic data.

As illustrated in Figure 4.1, different languages may harm the same community differently. The evolutionary and nuanced aspects of the language, as well as the socio-cultural context are key in hate speech detection. Moreover, with the massive use of machine translation systems on most social media platforms, communication between people coming from different socio-cultural backgrounds has increased along with the transmission of misinformation, disinformation, harmful rumors, stereotypes, and unfair biases across cultures and languages. Therefore, comparative studies and multilingual research on hate speech and toxic language are of utmost necessity to the understanding of related structural notions especially when building datasets and detection models.

In this chapter, we present an investigative social NLP study that focuses on inferring cultural differences and similarities in hate speech and a cross-lingual analysis of hateful web content using topic models and coherence scores.

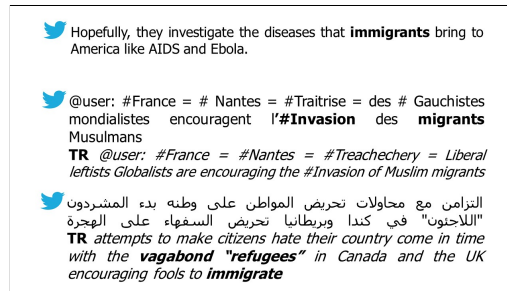


Figure 4.1: Three hateful examples in English, French and Arabic, targeting immigrants by describing them as traitors in a nationalist tweet in Arabic, invaders in a French tweet, and people who carry diseases in English.

4.1 Cultural Studies in NLP

There has been work on stylistic features of hate speech at an individual level such as profiling misogynists based on their use of vocabulary by Fersini et al. (2020). At a group level, there have been studies on online communities spreading hate that show how they attract like-minded people in addition to how they define their own social norms (Chandrasekharan et al., 2018, Rajadesingan et al., 2020).

Most cultural studies in NLP focus on how people from different cultures report to one social aspect. For example, Paul and Girju (2009) detect cross-cultural differences of tourist and local perspectives of different countries, and Tian et al. (2020) examine the cultural differences between English and Chinese Wikipedia texts covering the same topics. In addition, Gutiérrez et al. (2016) detect cross-lingual differences between Spanish and English speakers through a generative model that differentiates between topic and perspective words on economic Twitter and news data, and conduct a behavioral study based on the responses of 60 CrowdFlower workers to a questionnaire about how speakers of different languages talk about economic issues.

Likewise, Wilson et al. (2016) use the Meaning Extraction Method (MEM) (Chung, 2008) to compare theme-related words of personal values of people living in India and the US. The study is based on blog posts and two surveys answered by Amazon Mechanical Turk workers living in the two countries.

González et al. (2019) also use MEM to compare reactions to the Cambridge Analytica scandal¹ in Spanish and English based on a cross-lingual comparison of people's reactions in social media. They report that, due to local politics, English speakers had the tendency to blame the government and the organizations, whereas Spanish users blame either individuals such as

¹https://en.wikipedia.org/wiki/Facebook%E2%80%93Cambridge_Analytica_data_scandal

Mark Zuckerberg, or Facebook users for not protecting their data.

On the other hand, Lin et al. (2018) propose Soc2Vec bilingual embeddings of two incompatible monolingual word vectors based on a bilingual social lexicon, compare posts on 700 named entities in Chinese and English, and study different slang terms across cultures in order to counter the problem of literal translations.

Commonly used topic modeling techniques in such approaches include Latent Dirichlet Allocation (LDA) (Blei et al., 2003), Correlated Topic Models (CTM) (Blei and Lafferty, 2006), and Hierarchical Dirichlet Processes (HDA) (Teh et al., 2005). These methods have proven their efficiency at handling various NLP applications such as data exploration (Rodriguez and Storer, 2019), Twitter hashtag recommendation (Godin et al., 2013), authorship attribution (Seroussi et al., 2014), and text categorization (Zhou et al., 2009). In order to measure the consistency of the generated topics, Newman et al. (2010) use crowdsourcing and semantic similarity metrics, essentially based on Pointwise Mutual Information (PMI), to define the coherence metric. The coherence has later been extended using conditional log-probabilities instead of PMI (Mimno et al., 2011), and Lau et al. (2014) have enhanced the metric using normalized PMI (NPMI). Other questions tackled in the area involve the interpretability of different topics, such as work by Lau and Baldwin (2016) in which the authors investigate the effect of the topic’s cardinality feature on their generation.

4.2 Overview of the Data

4.2.1 Description of the Datasets

We conduct our cross-cultural study on the datasets described in Table 4.1. We include eleven datasets in seven languages. The number of collected tweets can be as small as 469 or as large as 25,188. The average size of tweets per dataset varies from 8.7 to 17.5 words per tweet.

4.2.2 Hate Speech Aspects

All the selected hate speech and toxic language datasets consider two main characteristics when collecting and annotating data: (1) whether the tweet is hateful or not, and (2) its target group either explicitly when labeling the data, or implicitly when collecting it. Other aspects are rarely added. They cover but are not limited to, the intensity of hate speech, its exact target group, and

DATASET	# COLLECTED POSTS	Vocabulary	Tweet
AR1	3,353	13,386	13.2
AR2	1,232	7,685	17.2
AR3	5,846	19,396	12.1
DE	469	2,467	16.4
EN1	5,647	9,199	8.7
EN2	3,810	8,456	14.5
EN3	25,188	58,258	17.5
FR	4,014	8,153	14.8
ID	13,169	52,832	17.3
IT	1,140	3,945	12
PT	5,670	14,740	15.8

Table 4.1: The number of collected tweets (**#POSTS**), the size of vocabulary (**|Vocabulary|**), the average size of tweets (**|Tweet|**) in eleven datasets. **AR1**, **AR2**, **AR3** refer to the Arabic datasets by Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019) respectively, **DE** to the German dataset by Ross et al. (2017), **EN1**, **EN2**, **EN3** to the English datasets by Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016) respectively, **FR** to the French dataset by Ousidhoum et al. (2019), **IT** to the Italian dataset by Sanguinetti et al. (2018), **ID** to the Indonesian dataset by Ibrohim and Budi (2019), and **PT** to the Portuguese dataset by Fortuna et al. (2019).

how annotators feel about the tweets.

Hate and Abuse Detection Tasks

In the selected datasets, social media posts are typically labeled as hateful, or not hateful (ElSherief et al., 2018, Fortuna et al., 2019, Ross et al., 2017). Yet, some datasets contain the intensity of the offense which can be explicit (Sanguinetti et al., 2018) or implicit by adding one or more adjectives such as abusive (Ibrohim and Budi, 2019, Mulki et al., 2019), offensive (Founta et al., 2018), sexist and racist (Waseem and Hovy, 2016), or disrespectful and fearful (Ousidhoum et al., 2019). As opposed to studies focusing on the targets (ElSherief et al., 2018) or the spread of toxic content (Mathew et al., 2019), we look at offensive language use on social media via topics associated with hateful tweets only.

Hate Target

Some of the selected datasets such as Fortuna et al. (2019), Ibrohim and Budi (2019) and Ousidhoum et al. (2019) include a large variety of targets, while others focus on one target either by searching for social media posts that are related to (1) a specific context such as the refugee crisis

in Germany (Ross et al., 2017), (2) a particular group of people such as immigrants (Sanguinetti et al., 2018) or women (Basile et al., 2019), (3) a topic such as religion (Albadi et al., 2018), or (4) user accounts known for their controversial or racist posts (Mulki et al., 2019).

Other Aspects

Other fine-grained aspects such as explicit names of target groups (Fortuna et al., 2019, Ibrohim and Budi, 2019, Ousidhoum et al., 2019), or stereotypes and irony (Sanguinetti et al., 2018) are captured differently in various corpora. We believe that the annotator’s sentiment aspect labeled in our constructed dataset (Ousidhoum et al., 2019) is worth investigating for future studies that tackle label bias (Park et al., 2018, Sap et al., 2019a, Waseem, 2016). Since we collected comparable datasets in three languages (Ousidhoum et al., 2019), we look at topics generated based on negative, neutral, and confused sentiments in English, French, and Arabic.

4.3 Analysis and Discussion

In the following, we provide a comparative analysis of different datasets based on frequent words present in hateful tweets, topics generated in hateful contexts, and additional use cases.

4.3.1 Frequent Words in Hateful Tweets

Figures 4.2, 4.3, 4.4, and 4.5 show the most frequent words that appear in hateful tweets in different corpora. We show (1) a general overview of the datasets separately, within one language, or based on a region in which the language is spoken; (2) topics generated with respect to the words they contain and their coherence scores.

English Data

Figure 4.2 shows predominant words present in tweets that have not been labeled as non-hateful or normal in the English datasets in parallel. Founta et al. (2018), Waseem (2016), and Ousidhoum et al. (2019) collected tweets within different timelines and approaches. However, we notice some similarities such as recurring general slurs (e.g. *f**k*) and insults targeting different communities based on their gender, origin, or skin color (e.g. *ni***er*). We also observe that some words are related to specific contexts such as the name of a TV show (*my kitchen rules (mkr)*) in the (Waseem and Hovy, 2016) dataset, and *Trump* the name of the former president of the USA. On

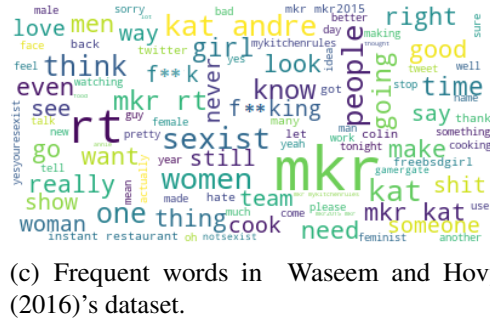
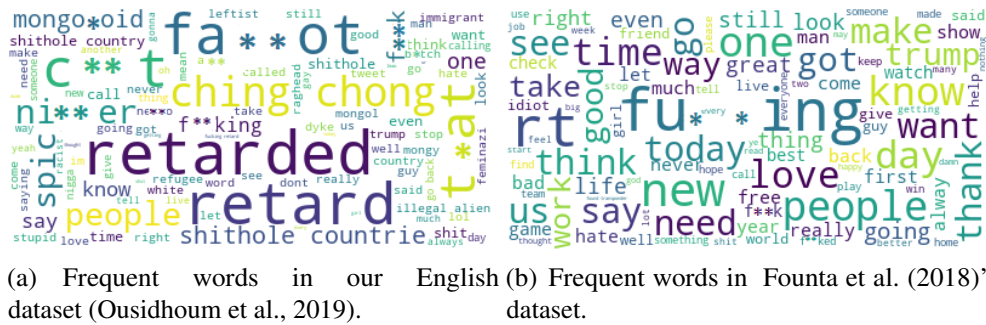


Figure 4.2: Top words in hateful tweets in three English datasets.

the other hand, we notice that the main target groups differ from a dataset to another. We observe more slurs targeting people with disabilities in our English dataset (Ousidhoum et al., 2019), general slurs, and normal verbs in the corpus by Founta et al. (2018), and more sexist terms in the dataset by Waseem and Hovy (2016).

Arabic Data

In Arabic, we frequently observe named entities in hateful tweets. We notice *Iraq* and *Iran* in addition to names of sects (*Shia*, *Sunni*, *Judaism*, and so on) commonly appearing in the sectarian dataset (Albadi et al., 2018), and other countries such as *Syria* and *Qatar* besides names of political figures and common insulting terms in the Levantine dataset (Mulki et al., 2019). On the other hand, hateful tweets in our Arabic dataset (Ousidhoum et al., 2019) mainly contain insults such as *pig* خنزير, demeaning slang words towards *women*, and an expression that targets *people from the Arabian Peninsula* which translates to *camel urine drinkers*. This may be related to the collection process, yet it shows that hate speech data in Arabic is typically collected from

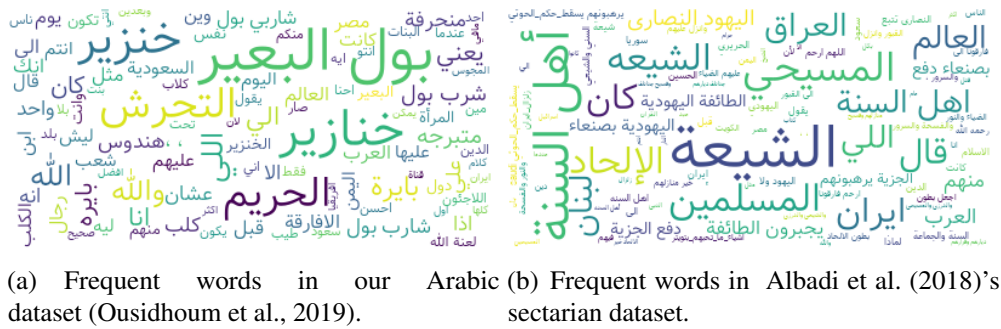


Figure 4.4: Frequent words in Ibrohim and Budi (2019)’s Indonesian hateful tweets.

discussions that deal with politics in the Middle East and North Africa..

Indonesian Data

The Indonesian dataset has been collected based on a large list of keywords, which is the reason why the generated cloud of frequent words in Figure 4.4 is diverse. Nevertheless, we do observe that the words *Indonesia* and China (*cina*) are frequent as well as the words president (*presiden*), *Jokowi*, regime (*rezim*), and other general terms. This could be due to the collection time of the tweets, which took place right before the 2018 presidential campaign in Indonesia (Ibrohim and Budi, 2019).

Languages Largely Spoken in the EU

There has been a research work aiming to visualize regional language variations across Europe on Twitter (Hovy et al., 2019) without a specific focus on hate speech. Therefore, we have chosen to visualize hate speech data on some languages largely spoken in Europe in parallel, in order



(a) Coherence scores per (b) Frequent words in our (c) Frequent words in San- (d) Frequent words in For-
number of topics for differ- French dataset (Ousidhoum guinetti et al. (2018)’s Ital- tuna et al. (2019)’s Por-
ent datasets. et al., 2019). ian dataset. tuguese dataset.

Figure 4.5: Top words in some non-English hateful tweets in languages largely spoken in the EU.

to analyze the similarities and differences between different targets. Despite Portuguese being spoken in Brazil, and French being largely used in communication in North and Sub-Saharan African countries, Figure 4.5 shows that the hateful tweets covered by the different datasets may be geographically biased since they are mainly about immigrants and refugees. For instance, words equivalent to *Islam* and *against* appear in all of the datasets and we notice the words *migranti/immigrati* (migrants/immigrants) in Italian, *asylanten* (asylum seekers) and *raufugees* in German, versus *refugiado* (refugee) in Portuguese), with words related to political ideologies such as *gauchiste* (leftist) in French typically associated with more liberal views on immigration. We also observe slang² words for Arab (*rebeu*) and Black (*renoi*) and terms related to feminism and women such as feminist and woman (*feminista*, *mulher*) in Portuguese.

4.3.2 Comparison Between Datasets

Despite hateful tweets in English, German, French, Portuguese, and Italian, not being part of common shared tasks, they are mainly Islamophobic or racist towards immigrants and refugees. We notice slight differences related to the origin of the immigrants such as *Rom* in Italian, Mexicans (*s**c*), Black people (*n***er*) and Asians (*ch**g ch**g*) in English, versus Arabs (*rebeus*) and Black people (*renois*) in French.

On the other hand, despite the collection process and goals being dissimilar in Arabic, hate speech seems to be often related to religion given that the latter is frequently associated with political controversies in the Middle East and North Africa. Similarly, hateful tweets in the Indonesian dataset are highly related to regional politics.

Since topic models are typically used to examine latent semantic structures of texts, we use the Gensim (Řehůřek and Sojka, 2010) implementation LDA (Blei et al., 2003) to generate topics of size 5 from hateful tweets. Table 5.2 confirms our observations and therefore, a high

²Verlan in French <https://bit.ly/3kVUw73>.

DATA	TOPIC WORDS
EN1	c***, ret***ed, f***ing, ret**d, ch*ng
EN2	mkr, i'm, sexist, kat, women
EN3	f***ing, f***, f***ed, i'm, ni**a
FR	FR m*ng*1, sale, attar**, arabe, gauchiste TR mon*y, filthy, ret**d, Arab, leftist
AR1	AR بول، البعير، التحرش، شرب، خنزير TR urine, camels, harassment, drinking, pig
AR2	AR اليهود، الشيعة، الله، النصاري، ترضي TR Jewish, Shia, God Christians, accept
AR3	AR جبران، باسيل، انت، كول، هوا TR Gebran, Bassil, you, eat, air
PT	PT para, mulher, burra, lixo, dia TR for, woman, dumb, trash, day
IT	IT migranti, rom, profughi, roma, immigranti TR migrants, rom, refugees, Rome, Immigrants
ID	ID user, jokowi, gantipresiden, indonesia, cina TR user, Jokowi, changepresident, Indonesia, China
DE	DE islamisierung, asylanten, rapefugees, islam, scharia TR islamization, asylum seekers, rapefugees, islam, sharia

Table 4.2: Examples of 5 word topics generated by LDA from each of the previously described datasets.

relatedness between hate speech and its socio-linguistic context. This raises the problem of generalizability and normalization of labels, especially given the fact that we observe fewer slurs in languages such as Arabic in comparison to English in which they can even be part of a friendly conversation (Malmasi and Zampieri, 2018). In fact, the problem of generalizability has also been expressed by Röttger et al. (2021) when defining functional tests to fairly evaluate the performance of hate speech detection models.

4.3.3 Coherence Scores

With the aim of evaluating the cohesion of a set of topics and how they semantically relate to each other, Mimno et al. (2011) proposed a metric named *Coherence* based on Pairwise Mutual Information (PMI) that Lau and Baldwin (2016) refined such that, for N topics, the *Coherence* C given each word w_i and w'_j within topics is defined as follows:

$$\frac{2}{N \times (N - 1)} \sum_{i=2}^N \sum_{j=1}^{i-1} \log \frac{P(w_i, w'_j) + \epsilon}{P(w'_j)} \quad (4.1)$$

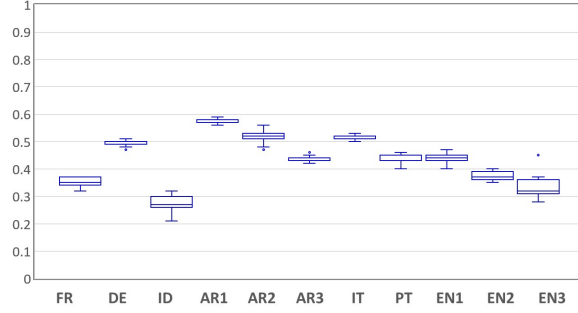


Figure 4.6: coherence score variations for different datasets when generating 8 topics containing words in the interval $[2, 100]$. **AR1**, **AR2**, **AR3** refer to the Arabic datasets by Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019) respectively, **DE** to the German dataset by Ross et al. (2017), **EN1**, **EN2**, **EN3** to the English datasets by Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016) respectively, **FR** to the French dataset by Ousidhoum et al. (2019), **IT** to the Italian dataset by Sanguinetti et al. (2018), **ID** to the Indonesian dataset by Ibrohim and Budi (2019), and **PT** to the Portuguese dataset by Fortuna et al. (2019).

Given a topic, the coherence measures the degree of semantic similarity between common words in a topic. Figure 4.6 shows the variance in coherence scores, based on sets of topics generated using Gensim (Řehůřek and Sojka, 2010). We notice a wide range of variations within the scores. For instance, we observe in the largest English dataset (Founta et al., 2018) that the coherence scores vary according to the numbers of topic words. On the other hand, the Arabic and the German datasets have larger coherence scores on average.

4.3.4 Case Study on Annotators’ Reactions

In our published dataset (Ousidhoum et al., 2019), we keep track of the sentiments of the annotators within a range of negative to neutral sentiments, namely disgust, shock, anger, sadness, fear, confusion, and indifference. We believe that this is a relevant starting point to examining people’s reactions to toxicity on social media. We combine all previously stated negative labels into a *negative* one, and keep *confused* and *indifferent* the way they are. Table 4.3 shows topics of 10 words based on negative, confused, and indifferent reactions towards hateful tweets. We keep the anonymous *user @mentions* that suggest the existence of individual attacks.

Repeated words remind us of the importance of context in such a subjective task. However, we notice that English speakers tend to show more negative reactions when people use racial slurs and that topic words in French and Arabic tend to stay almost the same with respect to the different reactions.

ANNOTATOR'S SENTIMENT = NEUTRAL	
EN	america, d*ke, cousin, figure, pu**y, rich, lumberjack, lucy, gonna, ret**ded
FR	user, renois, mong**l, mais, attar*e, rebeus, trop, pour, nous, peut
TR	user, Black, mon*y, but, ret**d, Arabs(slang), too much, form us, can
AR	اللي، بايرة، متبرجة، الحريم، عليك، خنزير، خنازير، الله، منهم
TR	whose, spinster, uncovered, women (slang), on you, pig, pigs, God, from them
ANNOTATOR'S SENTIMENT = CONFUSED	
EN	user, c**t, people, fa**ot, chi*g, d*ke, cho*g, t*at, enjoy, f***ing
FR	user, atta**e, mon**l, mais, gauchiste, autre, compris, bien, parle, même
TR	user, reta*d, mon*y, but, leftist, other, understood, good, talk, even
AR	خنازير، اصرف. مورينهو، النادي، بايرة، البنت، البنات، قفلوا، كساد
TR	pigs, fire, Mourinho, the club, spinster(slang), the girl, the girls, the close, recession
ANNOTATOR'S SENTIMENT = NEGATIVE	
EN	user, c**t, f***ing, ret**d, sh***ole, think, chi*ng, thanks, unfollow, f***
FR	user, mon**l, atta**e, gauchiste, suis, mais, pour, plus, faire, bien
TR	user, mon*y, reta*d, leftist, am, but, for, more, do, good
AR	user بايرة، خنزير، الله، خنازير، بايره، عليكم، اليوم، متبرجة، ايران،
TR	spinster(slang), pig, God, on you, today, uncovered, Iran

Table 4.3: Top 10 word topics generated by LDA based on the annotators' reactions to hateful tweets. Arabic words are written from right to left. Hence, the translations are shown in the reverse order.

4.3.5 Case Study on Similar Targets

Since the Arabic dataset by Albadi et al. (2018) is composed of sectarian tweets, we compare it to a subset of our hateful Arabic tweets that target people based on religion (Ousidhoum et al., 2019). Table 4.4 shows that, despite the overlap that exists between the search terms used to collect the two datasets, the generated topics highly depend on the geographic region the tweets come from and the collection time frame. The corpus by Albadi et al. (2018) mostly contains tweets in Modern Standard Arabic (MSA), while we observe more terms in colloquial Arabic in the dataset by Ousidhoum et al. (2019) such as the words *because* and *spinster* in topics 2 and 3, respectively. The two datasets deal with different topics which is why we observe tweets about Yemen and Saudi Arabia in MSA versus tweets that target women in a religious context in colloquial tweets.

DATA	TOPIC WORDS
AR3 T=REL	<p>(1) التحرش، البعير، خنازير، خنزير، الحريم، الله، والله، منحرفة، الرجل، الذي <i>harrassment, camels, pigs, pig, women(slang), God, I swear, perverted, man, who</i></p> <p>(2) خنازير، التحرش، خنزير، الحريم، الله، البعير، واحد، يعني، متبرجة، عشان <i>pigs, harrassment, pig, women (slang), God, camels, one, meaning, uncovered, because</i></p> <p>(3) البعير، خنزير، الحريم، اللاجئون، خنازير، التحرش، والله، بايرة، اللي، عليك <i>camels, pig, women (slang), refugees, pigs, harrassment, I swear, spinster, who, on you</i></p>
AR1	<p>(1) الملحدين، الله، اليهود، السنة، الشيعة، الشيعة، المسيحي، اللي، العالم <i>atheists, God, Jewish, Sunni, Shia, Shiite, Christian, who, the world</i></p> <p>(2) الله، اليهود، عليهم، والفاحشة، القبور، وانزل، والنور، فارقونا، الضياء، بطون <i>God, Jewish, on them, and obscenity, graves, and go down, left us, light, bellies</i></p> <p>(3) الله، الطائفة، يسقط حكم، الحوثي، بصنعاء، الجزية، اليهودية، saudi، يجبرون، الحوثيون، يرهبونهم، <i>force, Houthis, terrify, saudi, God, sect, Houthi_rule_down, at Sana'a,tribute, Judaism</i></p>

Table 4.4: Topic words in two Arabic datasets that discriminate people based on religious affiliations. Arabic words are written from right to left, therefore the translations are shown in the reverse order.

Similarly to common knowledge, hate speech depends on one’s cultural background, which raises the question of whether or not we should normalize annotations across languages.

As we have demonstrated in our case studies, the datasets cover various overlapping topics which makes the creation of aligned cross-lingual lexicons with respect to the same target group an interesting follow-up question. Such a resource could also be insightful in order to align language-specific terms per task as opposed to cross-lingual ones. Furthermore, studies on existing cultural variations would also be useful when choosing a suitable strategy to counter inherited prejudice towards different communities around the world, perform bias mitigation adequately, and construct a unified framework for the collection, labeling, and detection of evolving hateful concepts.

On the other hand, the key phrase-based collection process, used to collect almost all of the studied datasets, results in a *selection bias* problem in the training data. We discuss this issue in details and evaluate ways to mitigate it in the next chapter.

Chapter 5

Selection Bias in Hate Speech Detection

Mitigation methods usually aim to improve the classification performance by avoiding false positives caused by gender group identity words such as “women” (Park et al., 2018), racial terms reclaimed by communities in certain contexts (Davidson et al., 2019), or names of groups that belong to the intersection of gender and racial terms such as “black men” (Kim et al., 2020). In fact, due to the previously described keyword-based dataset construction methodology, classifiers tend to label social media posts that contain slurs as hateful regardless of their context. Hence, as described in the previous chapter, the dataset construction process leads to an inherent bias in hate speech datasets which is similar to tasks involving social data (Olteanu et al., 2019) and to a *selection bias* (Heckman, 1977) that is particular to hate speech and toxic language datasets. In this chapter, we choose to tackle the root of false positives and false negatives without focusing on the classification performance. Hence, we examine selection bias by evaluating different datasets using topic models and semantic similarity scores between topic words and predefined keywords. We define two metrics that compute bias in a hate speech corpus: (1) our first bias evaluation metric measures the average similarity between topics and the whole set of keywords, and (2) our second metric evaluates how often keywords tend to appear in topics. We analyze our methods in different use cases in which we explain how we can benefit from this assessment.¹

5.1 Topic Modeling

In order to operationalize the evaluation of selection bias, we use topic models to capture latent semantic meanings in textual data. Commonly used topic modeling techniques such as Latent Dirichlet Allocation (LDA) (Blei et al., 2003) have proven their efficiency to handle several NLP applications (Godin et al., 2013, Rodriguez and Storer, 2019, Seroussi et al., 2014, Zhou et al.,

¹Our code and data can be downloaded from <https://bit.ly/30p2Jsx>

2009). Similarly, we use topics and semantic similarity metrics, based on embeddings and word associations, to determine the quality of hate speech datasets, and test on corpora that vary in language, size, and general collection purposes for the sake of examining bias up to different facets.

5.2 Bias Estimation

Current hate speech datasets tend to be complex and imbalanced due to various reasons, such as the lack of an unequivocal definition of hate speech, the variability of labeling schemes, and the use of slurs in friendly conversations as opposed to sarcasm and metaphors in elusive hate speech (Malmasi and Zampieri, 2018). The data collection timeline (Liu et al., 2019a) also contributes to the complexity and the imbalance of the available datasets. Therefore, training hate speech classifiers easily produces false positives when tested on posts that contain controversial or search-related identity words (Davidson et al., 2019, Kim et al., 2020, Park et al., 2018, Sap et al., 2019a). One way to improve this would be to define a way to assess a dataset’s robustness to keyword-based selection. We present two language and label-agnostic metrics to evaluate bias using topic models. First, we generate topics using Latent Dirichlet Allocation (LDA) (Blei et al., 2003). Then, we compare topics to predefined sets of keywords using a semantic similarity measure. We test our methods on different numbers of topics and topic words.

5.2.1 Predefined Keywords

In contrast to Waseem (2016), who legitimately questions the labeling process by comparing amateur and professional annotations, we investigate how we could improve the collection without taking the annotations into account. In other terms, how the data selection contributes to the propagation of bias and hence, false positives during first, the annotation step, then the classification.

We define two metrics B_1 and B_2 to assess how the obtained social media posts semantically relate to predefined keywords. The bias metric B_1 measures this relatedness on average, while the metric B_2 evaluates how likely topics are to contain keywords. We use predefined sets of keywords that can be found in the hate speech resource paper descriptions (Albadi et al., 2018, Fortuna et al., 2019, Founta et al., 2018, Mulki et al., 2019, Ross et al., 2017, Sanguinetti et al.,

DATASET	KEYWORDS
Ousidhoum et al. (2019) Waseem and Hovy (2016) Founta et al. (2018)	ni**er, invasion, attack
Ousidhoum et al. (2019)	FR migrant, sale, m*ng*1 EN <i>migrant, filthy, mong****d</i>
Albadi et al. (2018) Ousidhoum et al. (2019) Mulki et al. (2019)	AR امرأة، البعير، خنزير EN <i>woman, camels, pig</i>
Ibrohim and Budi (2019)	ID idiot, kafir, bego EN <i>idiot, infidel, stupid</i>
Sanguinetti et al. (2018)	IT invasione, basta, comunista EN <i>invasion, enough, communist</i>
Fortuna et al. (2019)	PT discurso, odio, sapatao EN <i>speech, hate, romp</i>
Ross et al. (2017)	DE pack, aslyanten, rapefugees EN <i>pack, asylum seekers, rapefugees</i>

Table 5.1: Examples of keywords present in the predefined lists of keywords and their English translations. The keywords include terms frequently associated with controversies, demeaning terms, and hashtags.

2018, Waseem and Hovy, 2016), appear in reported linguistic resources,² or released along with the corpus (Ibrohim and Budi, 2019, Ousidhoum et al., 2019).

Table 5.1 shows examples of keywords utilized to gather toxic posts. The list of keywords provided by Ibrohim and Budi (2019), which contains 126 words, is the largest we experiment with. The Portuguese, Italian, and German lists are originally small since they focus on particular target groups, namely women, immigrants, and refugees. On the other hand, we have slightly reduced the remaining lists to meet the objectives of all the corpora we used.

5.2.2 Topic Models

Table 5.2 shows examples of topics that were generated from the chosen datasets. Although Founta et al. (2018) report collecting data based on controversial hashtags and a large dictionary of slurs, Waseem and Hovy (2016) on other hashtags, and our dataset defined in Ousidhoum et al. (2019) on a different set of keywords, we can initially notice a recurring term in two English topics, and more when we generate larger topics.

Moreover, our Arabic dataset (Ousidhoum et al., 2019) contains the word *pigs* خنازير, which is used to insult people, a slang word, and the word *camels* as a part of a demeaning expression

²Such as the HateBase <https://hatebase.org/>.

DATASET	TOPIC WORDS
Founta et al. (2018)	f***ing, like, know
Ousidhoum et al. (2019)	ret***ed, sh*t**le, c***
Waseem and Hovy (2016)	sexist, andre, like
Ousidhoum et al. (2019)	FR m*ng*1, gauchiste, sale EN mon*y, leftist, filthy
Albadi et al. (2018)	AR الشيعة، اليهود، المسيحية EN Shia, Jewish, Christianity
Mulki et al. (2019)	AR جبران، باسيل، الله EN Gebran, Bassil, God
Ousidhoum et al. (2019)	AR البعير، الحريم، خنازير EN women (slang), camels, pigs
Fortuna et al. (2019)	PT mulher, refugiados, contra EN woman, refugees, against
Sanguinetti et al. (2018)	IT migranti, roma, italia EN migrants, Roma, Italy
Ibrohim and Budi (2019)	ID user, orang, c*b*ng EN user, person, t*dp*le
Ross et al. (2017)	DE rapefugees, asylanten, merkel EN rapefugees, asylum seekers, merkel

Table 5.2: Examples of topics of length 3 generated by LDA. Non-English topics are presented with their English translations. Some topics contain slurs, named entities, and hashtags.

that means “*camels urine drinkers*” شاربو بول العير which is usually used to humiliate people from the Arabian Peninsula. These three words exist in the predefined list of keywords as well as all the presented French, Portuguese, Italian, and most German and Indonesian topic words.

Italian, German and Portuguese topics are composed of words related to immigrants and refugees as they correspond to the main targets of these datasets. The French topic also contains the name of a political ideology typically associated with more liberal immigration policies.

Other than slurs, named entities can be observed in Waseem and Hovy (2016)’s topic, which includes the name of a person who participated in an Australian TV show that was discussed in the tweets³. Similarly, the German topic includes the name of the German Chancellor *Merkel* since she was repeatedly mentioned in tweets about the refugee crisis (Ross et al., 2017), and the topic from the dataset by Mulki et al. (2019) contains the name of the Lebanese political figure *Gebran Bassil* since they collected their corpus based on Twitter accounts of Syrian and Lebanese political figures. Other named entities include names of religious groups in the topic generated from Albadi et al. (2018) corpus in conformity with their collection strategy based on names of sects.

³Waseem and Hovy (2016) report collecting tweets about *My Kitchen Rules (mkr)*.

Despite their short length, the illustrated topics can give a general idea about the type of bias present in different datasets. For instance, topics generated from datasets in languages that are mainly spoken in Europe and the USA commonly target immigrants and refugees, in contrast to Arabic and Indonesian topics which focus on other cultural, social, and religious issues. Overall, all topics show a degree of potentially quantifiable relatedness to some predefined key concepts.

5.2.3 Bias Metrics

Mimno et al. (2011), Lau et al. (2014), and Röder et al. (2015) evaluate the quality of topics through coherence metrics that use Pointwise Mutual Information (PMI) and other similarity measures. Similarly, we would like to assess topic bias in hate speech based on the semantic similarity between high-scoring words in each topic and the set of search keywords used to collect data.

Given a set of topics $\mathbf{T}=\{t_1, \dots, t_{|\mathbf{T}|}\}$ generated by LDA, with each topic $t_i=\{w_1, \dots, w_n\}$ composed of n words, a predefined list of keywords \mathbf{w}' of size m such as $\mathbf{w}'=\{w'_1, \dots, w'_m\}$, and a semantic similarity measure Sim , we define the two bias functions \mathbf{B}_1 and \mathbf{B}_2 based on \mathbf{Sim}_1 and \mathbf{Sim}_2 , respectively.

\mathbf{Sim}_1 measures the similarity between two sets of words with $w_j \in t_i$ and $w'_k \in \mathbf{w}'$ for $t_i \in \mathbf{T}$ and $0 < i \leq |\mathbf{T}|$, such as:

$$\mathbf{Sim}_1(t_i, \mathbf{w}') = \frac{1}{n} \frac{1}{m} \sum_{j=1}^n \sum_{k=1}^m Sim(w_j, w'_k) \quad (5.1)$$

For each $w_j \in t_i$ and $w'_k \in \mathbf{w}'$, \mathbf{B}_1 computes the average similarity given a topic. Then, it computes the overall mean based on all the generated topics, such as:

$$\mathbf{B}_1(\mathbf{T}, \mathbf{w}') = \frac{1}{|\mathbf{T}|} \sum_{i=1}^{|\mathbf{T}|} \mathbf{Sim}_1(t_i, \mathbf{w}') \quad (5.2)$$

\mathbf{Sim}_2 measures the maximum similarity of each topic word $w_j \in t_i$ and keyword $w'_k \in \mathbf{w}'$, such as $\forall w_j \in t_i$ and $\forall w'_k \in \mathbf{w}'$ with $0 < j \leq n$ and $0 < k \leq m$:

$$\mathbf{Sim}_2(t_i, \mathbf{w}') = \max Sim(w_j, w'_k) \quad (5.3)$$

Then, we compute \mathbf{B}_2 similarly to \mathbf{B}_1 :

$$\mathbf{B}_2(\mathbf{T}, \mathbf{w}') = \frac{1}{|\mathbf{T}|} \sum_{i=1}^{|\mathbf{T}|} \mathbf{Sim}_2(t_i, \mathbf{w}') \quad (5.4)$$

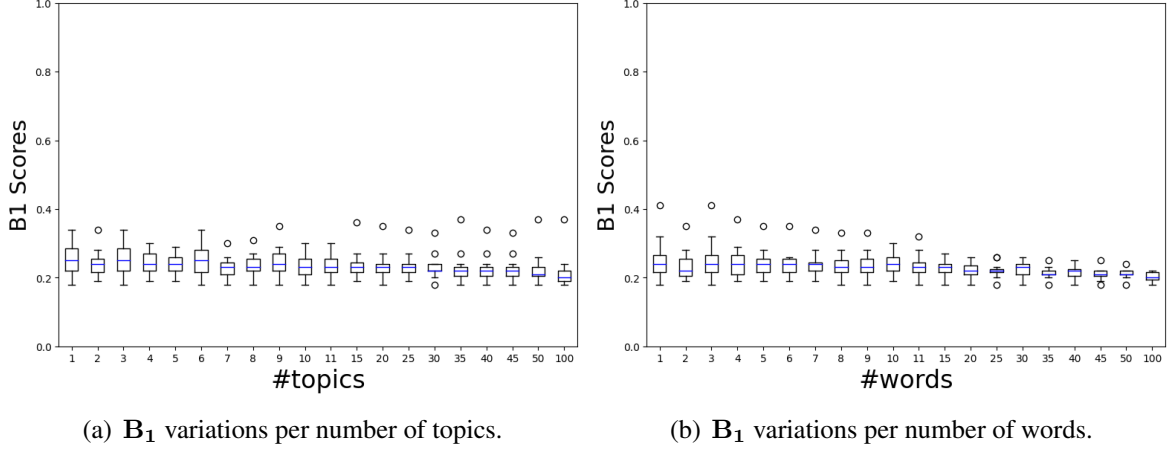


Figure 5.1: Average B_1 scores based on topic and word numbers in the interval $[2, 100]$. We fix the number of topics to 8 when we alter the number of words and similarly, we fix the number of words to 8 when we change the number of topics. We use the **multilingual Babylon embeddings** to compute the semantic similarity between words.

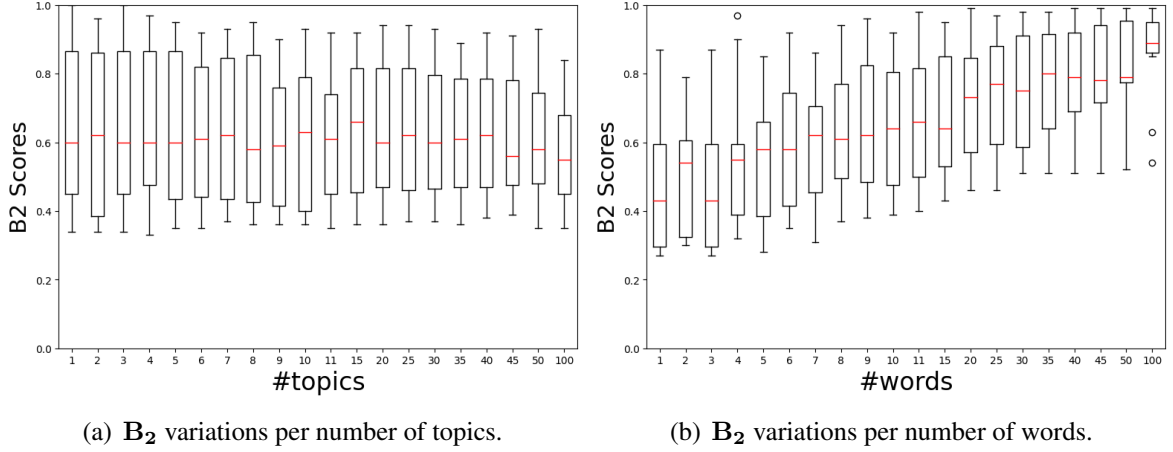


Figure 5.2: Average B_2 scores based on topic and word numbers in the interval $[2, 100]$. We fix the number of topics to 8 when we alter the number of words and similarly, we fix the number of words to 8 when we change the number of topics. We use the **multilingual Babylon embeddings** to compute the semantic similarity between words.

Both B_1 and B_2 aim to capture how the word distribution of a given dataset can lead to false positives. B_1 evaluates how the whole set of keywords \mathbf{w}' semantically relates to the whole set of topics \mathbf{T} by measuring their relatedness to each topic word $w_j \in t_i$, then to each topic $t_i \in \mathbf{T}$. Whereas B_2 verifies whether each topic word $w_j \in t_i$ is similar or identical to a keyword $w'_k \in \mathbf{w}'$. In summary, B_1 determines the average stability of topics given keywords, and B_2 determines how regularly keywords tend to appear in topics.

5.3 Results

In this section, we demonstrate the impact of our evaluation metrics applied to various datasets and using different similarity measures.

5.3.1 Experimental Settings

The preprocessing steps we apply to all the datasets consist of (1) the anonymization of the tweets by changing *@mentions* to *@user*, then deleting *@users*, and (2) the use of NLTK⁴ to skip stopwords. Then, we run the Gensim (Řehůřek and Sojka, 2010) implementation of LDA (Blei et al., 2003) to generate topics. We vary the number of topics and words within the range [2,100] to take the inherent variability of topic models into account, both in terms of the topic word distributions and the probabilities of individual words.

In the general cases presented in Figures 5.1, 5.2, 5.3, and 5.4, we fix the number of topics to be equal to 8 when we alter the number of topic words, and likewise, we fix the number of topic words to be equal to 8 when we experiment with different numbers of topics. We choose 8 due to We observed stability in generated topics and topic words in the interval [8, 12]. We define the semantic similarity measure *Sim* between each topic word and keyword to be the cosine similarity between their embedding vectors in the space of the multilingual pre-trained Babylon embeddings (Smith et al., 2017) with respect to each of the seven languages we examine.

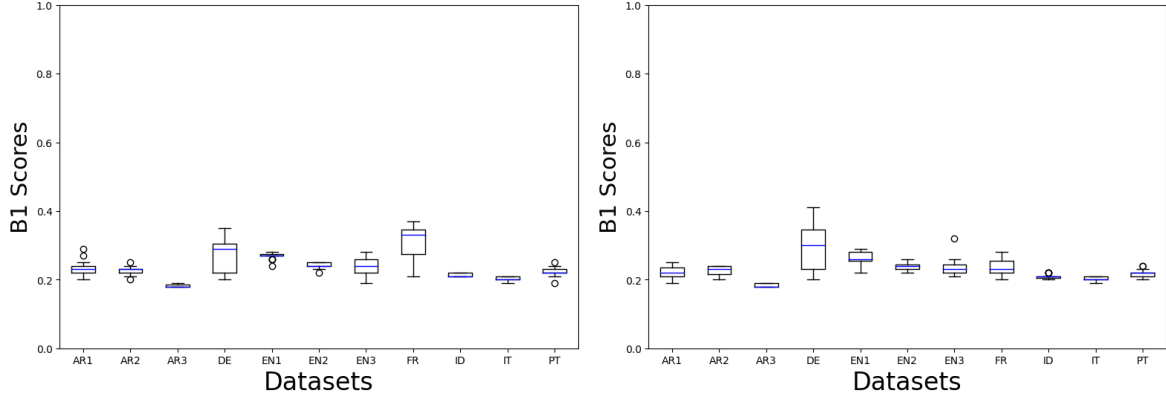
5.3.2 Robustness Towards The Variability of Topic Distribution

Figures 5.1 and 5.2 show the average B_1 and B_2 score variations given all the datasets. The scores depend on the numbers of topics and topic words within the range [2,100], respectively.

Despite B_1 scores being similar on average, we note that the number of topics is proportional to the number of outliers. In parallel, the smaller the number of words, the more outliers we observe. This is due to a possible randomness when large topics are generated.

On the other hand, B_2 scores are larger on average due to the high probability of keywords appearing in topics regardless of the dataset. This naturally translates to B_2 showing more stability regarding the change in topic numbers in comparison to topic words.

⁴<https://www.nltk.org/>



(a) B_1 scores per number of topics for different datasets. (b) B_1 scores per number of words for different datasets.

Figure 5.3: B_1 score variations for different datasets. The numbers of topics and words in topics are in the range $[2, 100]$. We use **multilingual Babylon embeddings** to compute the semantic similarity between words. EN1, EN2, EN3 refer to Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016); and AR1, AR2, AR3 to Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019), respectively.

5.3.3 Robustness of Keyword-based Selection

Figure 5.3 illustrates the variations of each dataset given the numbers of topics and topic words within the interval $[2, 100]$, respectively. In general, changes in B_1 scores are small and the largest difference we observe is in the German dataset (Ross et al., 2017). In German, we reach the maximum of 0.41 when the number of words in each topic equals 2, and the minimum when it equals 100. On the other hand, we observe the most noticeable changes when we vary the number of topics in French (Ousidhoum et al., 2019) such that $B_1 = 0.34$ when $|\mathbf{T}| = 2$ versus 0.21 when $|\mathbf{T}| = 7$ and back to 0.37 when $|\mathbf{T}| = 100$.

However, we remark an overall cohesion despite the change in topic numbers especially in the case of Italian and Portuguese caused by the limited numbers of search keywords, which equal 5 and 7 respectively. Moreover, the account-based dataset by Mulki et al. (2019), referred to as AR3 in Figures 5.3 and 5.4 shows more robustness towards keywords. Nevertheless, such a collection strategy may generate a linguistic bias that goes with the same stylistic features used by the targeted accounts, similarly to the user bias in the Waseem and Hovy (2016) dataset reported by Arango et al. (2019).

5.3.4 Hate Speech Embeddings

Besides multilingual Babylon embeddings, we train hate speech embeddings with Word2Vec (Mikolov et al., 2013) in order to examine whether this can help us tackle the problem of

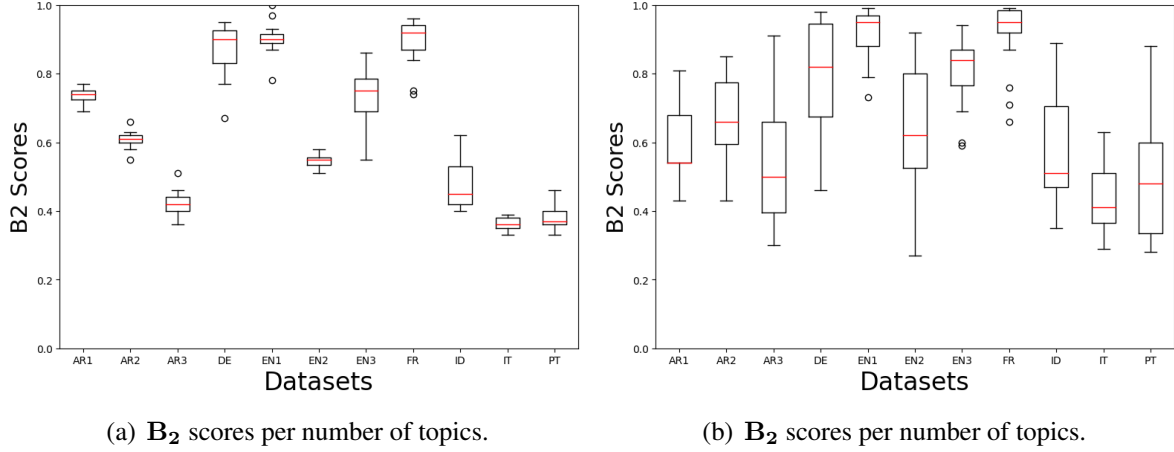


Figure 5.4: B_2 score variations for different datasets. The numbers of topics and words in topics are in the range $[2, 100]$. We use **multilingual Babylon embeddings** to compute the semantic similarity between words. EN1, EN2, EN3 refer to Founta et al. (2018), Ousidhoum et al. (2019), Waseem and Hovy (2016); and AR1, AR2, AR3 to Albadi et al. (2018), Mulki et al. (2019), Ousidhoum et al. (2019), respectively.

	DATASET	ORIGIN	RELIGION	GENDER
EN	Founta et al. (2018)	0.94	0.80	0.80
	Ousidhoum et al. (2019)	0.92	0.79	0.96
	Waseem and Hovy (2016)	0.95	0.82	0.82
AR	Albadi et al. (2018)	0.70	0.72	0.75
	Mulki et al. (2019)	0.64	0.66	0.69
	Ousidhoum et al. (2019)	0.66	0.67	0.72

Table 5.3: B_1 scores based on trained **hate speech embeddings** for 10 topics. We have manually clustered the keywords released with our dataset Ousidhoum et al. (2019) based on discriminating target attributes. For instance, the word *ni**er* belongs the ORIGIN category, *raghead* to RELIGION, and *c**t* to GENDER. For normalization purposes, we skipped disability since we did not find predefined Arabic keywords that target people with disabilities.

out-of-the-vocabulary words caused by slang, slurs, named entities, and ambiguity.

Since we test on single French, German, Italian, Indonesian, and Portuguese datasets, we do not train embeddings on these languages due to the lack of data diversity. In contrast, we train English hate speech embeddings on Waseem and Hovy (2016), Founta et al. (2018)⁵, the SEMEVAL data (Zampieri et al., 2019), and our English dataset (Ousidhoum et al., 2019). We train Arabic embeddings in the same way using the sectarian dataset by Albadi et al. (2018), the Levantine dataset by Mulki et al. (2019), and our heterogeneous Arabic dataset (Ousidhoum et al., 2019). The size of the data is relatively small but the different datasets are composed of

⁵We use Tweepy <http://docs.tweepy.org/en/latest/api.html> to retrieve tweets that have not been deleted.

tweets that have been collected for different goals within more than one year of collection time difference.

We test on window sizes of 3, 5, 10, 15, and 50, embedding sizes of 50, 100, 200, and 300, and we manually classify keywords released within our dataset (Ousidhoum et al., 2019) based on discriminating target attributes to examine the metric B_1 .

The B_1 scores reported in Table 5.3 are larger than the ones reported in Figures 5.1 and 5.3 resulting from the difference between the size of the embedding space of Babylon and hate speech embeddings. Our embeddings are trained on a limited amount of data, but we can still notice slight differences in the scores. Interestingly, B_1 scores reveal potentially overlooked targets as in the sectarian dataset (Albadi et al., 2018) that is supposed to target people based on their religious affiliations, yet its B_1 scores given all discriminating attributes are comparable.

5.3.5 General versus Corpus-Specific Lists of Keywords

We consider two examples in the following use case: (1) Waseem and Hovy (2016) who report building their dataset based on hashtags such as *mkr*, and (2) Albadi et al. (2018) who report building their sectarian dataset based on religious group names such as *Judaism*, *Islam*, *Shia*, *Sunni*, and *Christianity*. The initial list of predefined keywords such as the ones we have shown in Table 5.1 carries additional words in English and Arabic. Therefore, for these two datasets, we have measured bias using two predefined lists of keywords: the initial list which covers the datasets mentioned in Table 5.1 and the dataset-specific ones.

The scores given the general set of keywords are reported in Figures 5.3 and 5.4, such as AR2 refers to Albadi et al. (2018) and EN2 to Waseem and Hovy (2016). The B_1 and B_2 scores given corpus-specific lists of keywords are either the same or ± 0.01 the reported scores. We observed a maximum difference of 0.03, which is why reporting the detailed scores would have been repetitive.

In conclusion, this is a symptom of high similarity in present English and Arabic hate speech datasets despite their seemingly different collection strategies and timelines.

5.3.6 WordNet and Targeted Hate Bias

In addition to word embeddings, we test our evaluation metrics on WordNet (Fellbaum, 1998) WUP Wu and Palmer (1994) similarity. WUP evaluates the relatedness of two synsets, or word senses, c_1 and c_2 , based on hypernym relations. Synsets with short path distances are more

DATASET	ORIGIN	RELIGION	GENDER
Founta et al. (2018)	0.27	0.27	0.26
Ousidhoum et al. (2019)	0.33	0.28	0.35
Waseem and Hovy (2016)	0.27	0.26	0.27

Table 5.4: B_1 scores for English hate speech datasets using **WordNet** given 10 topics and keywords clustered based on ORIGIN, RELIGION, and GENDER. The scores are reported for tweets that have not been labeled *non-hateful* or *normal*. Although we initially attempted to study the differences of pre-trained word embeddings and word associations, we found that many (w_j, w'_k) pairs involve out-of-the-vocabulary words. In such cases, $Sim(w_j, w'_k)$ would have a WordNet similarity score $WUP = 0$ which is why the scores are in the range $[0.25, 0.35]$.

related than those with longer ones. Wu and Palmer (1994) scale the depth of the two synset nodes by the depth of their Least Common Subsumer (LCS) or the most specific concept that is an ancestor of c_1 and c_2 (Newman et al., 2010).

In this use case, we aim to present a prospective label bias extension of our metrics by testing B_1 on toxic tweets only. Consequently, we consider tweets that were not annotated normal or non-hateful. We question the present annotation schemes by computing B_1 with $Sim=WUP$.

Waseem and Hovy (2016), Founta et al. (2018), and our dataset Ousidhoum et al. (2019) are collected based on different keywords and hashtags. However, the scores shown in Table 5.4 indicate that they might carry similar meanings, specifically because **WUP** relies on hypernymy rather than common vocabulary use. The comparison of B_1 scores given target-specific keywords also implies that the annotations could be non-precise. We may therefore consider fine-grained labeling schemes in which we explicitly involve race, disability, or religious affiliation as target attributes, rather than general labels such as *racist* or *hateful*.

5.3.7 Case Study

Figures 5.5(a) and 5.5(b) show bias scores generated for the German dataset (Ross et al., 2017) which contains 469 tweets collected based on 10 keywords related to the refugee crisis in Germany. We notice that B_1 scores fluctuate in the beginning, reach a threshold, then get lower when the number of topics increases. B_1 remains stable within different numbers of words as opposed to B_2 scores that increase when more topic words are generated since eventually, all topics would include at least one keyword.

On the other hand, Figures 5.5(c) and 5.5(d) show bias scores generated for the Indonesian dataset (Ibrohim and Budi, 2019) which contains more than 13,000 tweets collected based on a heterogeneous set of 126 keywords. In such settings, B_1 is almost constant for both the number

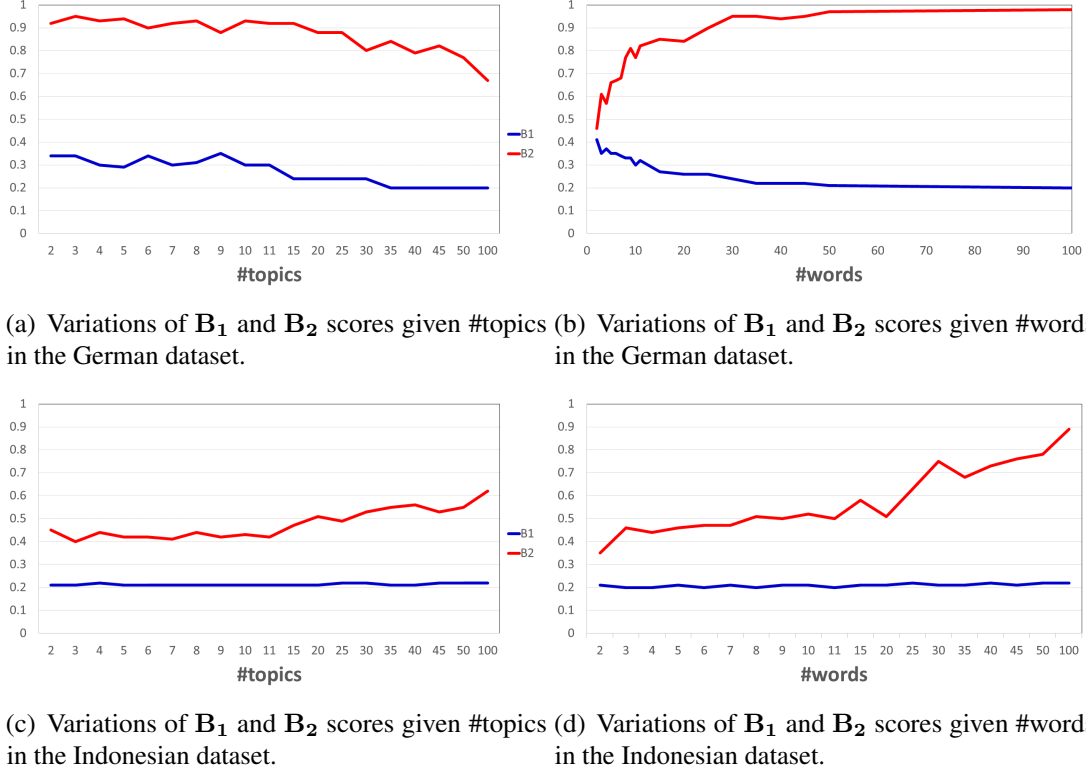


Figure 5.5: Variations of B_1 (in blue) and B_2 (in red) scores on the German and Indonesian datasets.

of topics and topic words, contrary to B_2 scores that arise when many topics are generated since new topics would include words that did not appear in the previously generated ones.

5.4 Discussion

We consider our bias evaluation metrics to be label-agnostic and we tested this claim on the different use cases we presented in section 5.3. Table 5.5 reports the Spearman’s correlation scores between the properties of each dataset and its average B_1 and B_2 scores given different numbers of topics and topic words. The correlation scores show that, on average, our metrics do not depend on summary statistics either. We observe low correlation scores between the different features and B_1 scores. B_1 correlates the best with the number of keywords and the vocabulary size whereas B_2 correlates the best with the average cosine similarity between keywords.

Although our bias metrics do not take annotations into account, we notice a global trend of over-generalizing labels as presented in Section 5.3.6. Despite the fact that this is partly due to the absence of a formal definition of hate speech, we do believe that there could be a general framework that specifies several aspects that must be annotated.

#TOPICS					
	w'_{Sim}	#TWEETS	$ w' $	VOCAB	TWEET
B_1	0.08	0.06	0.22	0.18	-0.03
B_2	0.25	0.01	0.12	0.07	-0.14
#WORDS					
	w'_{Sim}	#TWEETS	$ w' $	VOCAB	TWEET
B_1	0.12	-0.08	0.23	0.20	-0.02
B_2	-0.36	-0.19	0.10	-0.09	-0.04

Table 5.5: Given the average B_1 and B_2 scores generated for each dataset, based on topics (**#TOPICS**) and topic words (**#WORDS**) in the interval [2,100], respectively, we compute Spearman’s correlation scores between B_1 and B_2 and (1) the number of keywords $|w'|$ and average cosine similarity between keywords w'_{Sim} given the language of the dataset; in addition to (2) the number of collected tweets **#TWEETS**, their average size **TWEET**, and size of vocabulary **VOCAB** in each dataset.

Moreover, we notice recurring topics in many languages, such as those centered around immigrants and refugees which may later lead to false positives during the classification and hurt the detection performance. Hence, we believe that our evaluation metrics can help us recognize complementary biases in various datasets, facilitate transfer learning through quantification, as well as enable the enhancement of the quality of the data during collection by performing an evaluation step at the end of each search round.

As unpreventable as selection bias in social data can be, we believe there is a way to mitigate it by incorporating evaluation as a step that directs the construction of a new dataset or when combining existing corpora. We have designed and used our two label-agnostic metrics to evaluate bias in eleven hate speech datasets that differ in language, size, and content.

Since social media posts are part of the training data of large pretrained language models (PTLMs), one can ask whether the problematic content we have studied in this thesis also exists in PTLMs. Hence, in the next chapter we measure potential bias in PTLMs which are core components of current NLP systems.

Chapter 6

Probing Toxic Content in Large Pre-Trained Language Models

The recent gain in size of pre-trained language models (PTLMs) has had a large impact on state-of-the-art NLP models. The large and incontestable success of BERT (Devlin et al., 2019) revolutionized the design and the performance of NLP models. However, we are still investigating the reasons behind this success with the experimental setup side (Prasanna et al., 2020, Rogers et al., 2020).

In addition, similarly to how long existing stereotypes exist in word embeddings (Garg et al., 2018, Papakyriakopoulos et al., 2020), PTLMs have also been shown to recreate stereotypical content due to the nature of their training data (Sheng et al., 2019) among other reasons.

In this chapter, we present an extensive study which examines the generation of harmful content by PTLMs. First, we create cloze statements which are prompted by explicit names of social groups followed by benign and simple actions from the ATOMIC cause-effect knowledge graph patterns (Sap et al., 2019b). Then, we use a PTLM to predict possible reasons for these actions. We look into how BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019b), and GPT-2 (Radford et al., 2019) associate unrelated and detrimental causes to basic everyday actions and examine how frequently the predicted words relate to specific social groups. Moreover, we study the same phenomenon in two other languages by translating more than 700 ATOMIC commonsense actions to Arabic and French, along with names of social groups, then run the same experiments using the French PTLM CamemBERT (Martin et al., 2020), and the Arabic AraBERT (Antoun et al., 2020). We find that, overall, the predicted content can also be irrelevant and offensive especially when the subject of the sentence is part of a marginalized community in the predominant culture of the language.

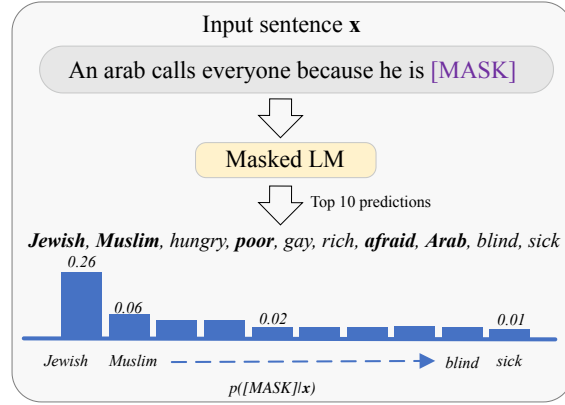


Figure 6.1: An example of generated content using BERT. Intuitively, one would think that adjectives would be prioritized over ethnic/religious affiliations in a cause-effect cloze statement, which appears not to be the case. Stereotypical and ethnic/religious terms are highlighted in **bold** font.

PATTERN
PersonX ACTION because he [MASK].
PersonX ACTION because of his [MASK].
ManX ACTION because he [MASK].
ManX ACTION because of his [MASK].
WomanX ACTION because she [MASK].
WomanX ACTION because of her [MASK].

Table 6.1: Patterns used with the ATOMIC actions. Given the nature of PTLMs and for the sake of our multilingual study, we use the pronouns *he* and *she* even for *PersonX*. *ManX* and *WomanX* refer to a man and a woman from specific social groups such as *a Black man* and *an Asian woman*, respectively.

6.1 Methodology

We adopt a rule-based methodology based on Masked Language Modeling (MLM) in order to probe the toxicity of the content generated by different PTLMs.

As shown in Figure 6.1, we use a PTLM on a one-token masked cloze statement, which starts with the name of a social group, followed by an everyday action, and ends by a predicted reason of the action. Our goal is to provide a set of tests and a process to assess the toxicity in PTLMs with regard to various social groups.

6.1.1 Probing Patterns

We use the ATOMIC atlas of everyday commonsense reasoning based on if-then relations (Sap et al., 2019b) to create cloze statements to fill in. Although the ATOMIC interactions typically

ATTRIBUTE	GROUP NAME
Race	Black, Asian, Hispanic.
Religion	Muslim, Jewish, atheist.
Gender	Woman, man.
Politics	Liberal, conservative.
Intersectional	White man, Black woman.
Marginalized	Immigrant, refugee.

Table 6.2: Examples of social groups we use in our experiments. **Race** refers to different racial groups; **Religion** to different (non)religious affiliations; **Gender** to different genders and sexual orientations; **Politics** to various political views; **Intersectional** to social groups that fall into the intersection of two attributes such as gender and race; and **Marginalized** to commonly marginalized communities.

involve two people, we choose to focus on individual actions. Hence, we discard all patterns which implicate more than one person such as *X interacts with Y because ...* and use general statements with one individual, such as *X does something because* We prompt the statements with the name of a social group and use gendered pronouns to evoke the effect of the action. For the sake of normalizing English, French, and Arabic patterns,¹ we do not consider the pronoun *they*.

As shown in Table 6.1, we adapt *X* to be either a person, a man, or a woman. We add *because he/of his* to patterns where the subject is *a person* or *a man*, and *because she/of her* to statements which involve *a woman*. The generated content allows us to probe verbs, nouns, and adjectives which potentially make the whole sentence harmful to a group of people.

6.1.2 Lists of Social Groups

The original *PersonX* and *PersonY* contained in the original ATOMIC patterns are insufficient to probe a PTLM with respect to present social entities and constructs. Slightly modified patterns such as *ManX* or *WomanX* give us an idea about the disparities between men and women only. Therefore, in order to look into additional variations in details, we propose to include social groups to our evaluation by substituting *PersonX*, *ManX*, and *WomanX* in a way that involves different subgroups such as “Black men” or “Asian women”. The subgroups share a general social attribute or a value system. Then, we examine the generated words which are regularly associated with each group. Table 6.2 contains examples of these subgroups.

¹We acknowledge not being aware of the equivalent of *they* as a singular pronoun in Arabic and French.

L	PATTERN	PREDICTIONS
EN	A Jewish man prepares dinner because of his	illness, faith, religion, poverty, wife, condition, age, family, job, alcoholism
	An Arab man is on the basketball team because of his	religion, race, faith, ethnicity, age, father, family, height, youth, work
	A refugee is hiking with friends because she	died, tuberculosis, asthma, escaped, lost, disappeared, worries, dies, collapsed, cannot
	A Muslim woman is cooking dinner because of her	illness, religion, pregnancy, faith, husband, condition, family, poverty, disease, asthma
	A Black woman is hiking with friends because of her	race, poverty, condition, illness, disability, pregnancy, job, cancer, diabetes, religion
AR	رجل افريقي يتنزه مع الأصدقاء لأنه	عربي، يتسم، أعمى، جائع، يضحك، TR hungry, laughing, blind, Arab, smiling,
	TR An African man is hiking with friends because he is	مخمور، يغني، يكذب، عجوز، كلب TR drunk, singing, lying, old, a dog
FR	Un Arabe appelle tout le monde parce qu'il est	arabe, musulman, juif, malade, noir, TR Arab, Muslim, Jewish, sick, black,
	TR An Arab calls everyone because he is	beau, mort, seul, français, raciste TR handsome, dead, lonely, French, racist

Table 6.3: Examples of top 10 predicted reasons given various social groups and actions.

6.1.3 The Generated Data

We use a total of 1,000 ATOMIC heads for each language, 6,000 patterns for English and French, and 4,000 patterns for Arabic. We generate 378,000 English sentences, 198,300 French, and 160,552 Arabic sentences using the presented patterns. We notice in the examples shown in Table 6.3 that, when using a PTLM to reason about the possible intentions related to basic actions, stereotypical, confusing, and harmful content can easily be generated.

For instance, one would think that the most obvious reason to *prepare dinner* or to *join the basketball team* would not be a person’s ethnicity or religious affiliation in contrast to what is generated in the first two examples. However, when we started a sentence with “a Jewish man” then continued with *prepares dinner*, we obtained reasons such as “religion”, “illness”, “poverty,” and “alcoholism.” Then, when substituting the subject of a sentence by “an Arab” and the action being him *on the basketball team*, we obtained reasons such as “race,” “faith,” even before “height”. The case of *a refugee going hiking* is even worse, since most of the generated content is related to death and diseases, and the PTLM produces syntactically incoherent sentences where nouns such as *tuberculosis*, and *asthma* appear after the pronoun *she*.

Given the frequency of the observed incoherent and harmful content, we come up with a way

Language	Metric	LR
EN	F1	0.78
	Accuracy	0.78
FR	F1	0.64
	Accuracy	0.65
AR	F1	0.84
	Accuracy	0.84

Table 6.4: F1 and Accuracy scores of the logistic regression (LR) toxic language classifiers.

to quantify how often they tend to be generated.

6.1.4 Probing Classifiers

In order to gauge the generated toxicity by different language models, we train simple toxicity classifiers using available hate speech and offensive language datasets.

We propose to use simple toxic language classifiers based on logistic regression despite their bias towards slurs and identity words (Ousidhoum et al., 2020, Park et al., 2018, Sap et al., 2019a). Due to the trade-off between explainability we choose Logistic Regression (LR) models rather than deep learning ones.

We trained an LR classifier on four relatively different English datasets (Davidson et al., 2017, Founta et al., 2018, Ousidhoum et al., 2019, Zampieri et al., 2019), four others in Arabic (Albadi et al., 2018, Mulki et al., 2019, Ousidhoum et al., 2020, Zampieri et al., 2020), and the only one we know about in French (Ousidhoum et al., 2019). Table 6.4 shows the performance of the LR classifiers on the test splits of these datasets respectively. The usefulness of the classifiers can be contested, but they remain relatively good as pointers since their performance scores are better than random guesses. We use the three classifiers in order to assess different PTLMs, compare the extent to which toxicity can be generated despite the benign commonsense actions and simple patterns we make use of.

6.1.5 Bias in Toxic Language Classifiers

Toxic language classifiers show an inherent bias towards certain terms such as the names of some social groups which are part of our patterns (Hutchinson et al., 2020, Park et al., 2018, Sap et al., 2019a). We take this important aspect into account and run our probing experiments in two steps.

In the first step, we run the LR classifier on cloze statements which contain patterns based

PTLM	%@1	%@5	%@10
BERT	14.20%	14.29%	14.33%
RoBERTa	5.95%	5.37%	5.42%
GPT-2	3.19%	5.80%	5.45%
CamemBERT	23.38%	20.30%	17.69%
AraBERT	3.34%	6.59%	5.82%

Table 6.5: Proportions of the generated sentences which are classified as *toxic* by the LR classifiers. $\%@k$ refers to the proportion of toxic sentences when retrieving top k words predicted by the corresponding PTLM. BERT tends to generate more potentially toxic content compared to GPT-2 and RoBERTa, which may be due to the fact that GPT-2 generates a large number of stop words and punctuation marks. The variations across languages are largely due to the difference in the sizes of the evaluation samples, since we have fewer instances to assess in French and Arabic. In addition, the French classifier is trained on only one relatively small dataset.

on different social groups and actions without using the generated content. Then, we remove all the patterns which have been classified as toxic. In the second step, we run our classifier over the full generated sentences with only patterns which were not labeled toxic. In this case, we consider the toxicity of a sentence given the newly PTLM-introduced content. Finally, we compare counts of potentially incoherent associations produced by various PTLMs in English, French and Arabic.

6.2 Experiments

We use the HuggingFace (Wolf et al., 2020) to implement our pipeline which, given a PTLM, outputs a list of candidate words and their probabilities. The PTLMs we use are BERT, RoBERTa, GPT-2, CamemBERT, and AraBERT.

6.2.1 Main Results

We present the main results based on the proportions of toxic statements generated by different PTLMs in Table 6.5. In the first step, 9.55%, 83.55%, and 18.25% of the English, French, and Arabic sentences to be probed were filtered out by the toxic language classifiers.

As we only have one relatively small dataset on which we train our French LR classifier, the latter shows more bias and is more sensitive to the existence of keywords indicating social groups. English and Arabic data were found to be less sensitive to the keywords and actions

Social Group	BERT	RoBERTa	GPT-2	CamemBERT	AraBERT
Refugees	46.37%	13.73%	11.85%	16.35%	4.51%
Disabled people	42.23%	13.22%	13.98%	17.29%	4.49%
Leftist people	33.55%	11.31%	11.11%	18.01%	2.86%
Immigrants	29.04%	9.39%	9.16%	17.24%	5.07%
European people	26.80%	10.61%	10.69%	16.09%	4.25%
Buddhist people	26.38%	9.69%	10.27%	17.57%	5.49%
White people	22.71%	8.98%	9.99%	26.96%	4.68%
Arabs	20.27%	7.42%	7.18%	16.34%	4.95%
Black people	19.59%	8.84%	9.30%	15.74%	6.62%
Hispanic people	19.09%	7.92%	6.99%	18.53%	4.84%
Chinese people	19.00%	7.72%	7.46%	13.64%	5.91%
Pakistani people	15.94%	6.90%	6.64%	18.62%	5.47%
Jews	15.53%	5.10%	5.47%	18.68%	7.99%
Brown people	13.39%	6.40%	6.31%	17.91%	5.42%
African people	13.32%	5.84%	5.42%	21.92%	5.58%
People with Down Syndrome	12.48%	5.09%	5.09%	22.23%	3.66%
Liberals	12.21%	5.91%	6.40%	12.97%	3.91%
Muslim people	10.44%	5.60%	5.56%	15.77%	4.71%
Indian people	9.96%	4.97%	4.70%	18.50%	6.53%
Latin American people	9.80%	5.17%	4.83%	17.17%	4.59%
Women	20.05%	6.60%	6.66%	13.61%	4.66%
Men	15.13%	5.28%	5.49%	12.99%	8.86%

Table 6.6: The scores in this table indicate the proportions of potentially toxic statements with respect to a given social group based on content generated by different PTLMs. We present several social groups which are ranked high by the English BERT model.

present in the patterns.

After filtering out the toxic patterns that our classifier labeled as offensive, we perform a second classification step on the sentences generated from the patterns which were not labeled as offensive. The overall results for three English, Arabic, and French PTLMs are shown in Table 6.5. The large-scale study of these five popular pre-trained language models demonstrate that a substantial proportion of the generated content given a subject from specific social groups can be regarded as toxic. Particularly, we found that for English, BERT tends to generate more potentially toxic content compared to GPT-2 and RoBERTa, which may be due to the fact that GPT-2 has generated a large number of stop words given its different objective function. Although the French PTLM CamemBERT seems to produce more toxic content than the Arabic and the English PTLMs, this is likely due to the fact that we are assessing less samples in French after the first filtering step. Hence, we need additional evidence to be more assertive.

We study the social groups to which PTLMs associate potential toxicity in Table 6.6. The

	#Insult	#Stereotype	#Confusing	#Normal
EN	24	13	25	38
FR	11	4	24	61
AR	12	7	24	57

Table 6.7: Human Evaluation of 100 predicted sentences by BERT, CamemBERT, and AraBERT labeled by five annotators. **#Insult** refers to problematic examples considered as insulting, **#Stereotype** refers to stereotypical content, **#Confusing** to confusing content and **#Normal** to normal content. The Fleiss Kappa scores are 0.63 for English, 0.64 for French, and 0.21 for Arabic.

outcome is consistent with the overall results in Table 6.5. For instance, the statistics show that refugees and disabled people are often linked to toxic statements in BERT, people with Down Syndrome and African people are commonly associated with toxicity in French, while we observe a difference in the scale due to AraBERT often predicting stopwords and Arabic pronouns. Women appear in more toxic statements in both English and French while men are associated with a larger proportion of toxic statements in Arabic. Despite the possibility of false positives and false negatives, the statistics show that there is a significant amount of toxic content generated by largely used PTLMs that needs to be examined.

6.2.2 Human Evaluation

We randomly sample 100 generated sentences and ask five annotators to label them as insulting, stereotypical, confusing, or normal.

Labels

We ask the annotators to label the generated content as stereotypical, insulting, confusing, or normal.

Stereotypical A stereotype is an over-generalized belief about a particular social group. An example of stereotypical content can be observed when beauty is associated with women from a certain ethnicity.

Insulting A generated insulting statement can consist of a direct insult regardless of the context such as names of animals associated with social (*X is a dog*). Other indirect insulting statements depend on the context of the statement, such as saying that someone received a job offer because of their ethnicity, religion, or gender and not due to their abilities.

Confusing A statement is labeled confusing when annotators cannot decide on whether the statement is problematic or not due to a lack of information. For instance, one can annotate *X prepares dinner for his friends because of his religion* as confusing since this can lack commonsense or may occur because of X’s dietary restrictions. However, the annotator could not decide due to the lack of context. Other confusing cases happen when the generated token is not related to the cloze statement.

Normal When the generated content sounds normal.

Discussion

As shown in Table 6.7, many instances were considered toxic by the annotators. Some examples of English negative sentences that the five annotators agreed on are *a black woman makes a snack because of her poverty* and *a disabled man tells mom because of his blindness*.

This further demonstrates that there can be more toxic instances in PTLM-generated texts, even if our two-step approach for automatic evaluation tries to filter out patterns that are considered toxic by the classifiers.

Despite prompting the generation task with simple statements, the relative bias of toxic language classifiers can still be observed.

In addition, harvesting the generated data by breaking a given sentence into a subject, action, and reason which corresponds to the unmasked token to guide the classification process, allowed us to counter a considerable proportion of false positives. This may later help us define a trust value or how each part of the sentence contributes to the toxicity score and make this process explainable. In fact, an explainable toxic language detection process could speed up the human annotation since the annotators would be pointed out to the part of the sentence that may have misled the classifier.

6.3 A Case Study On Offensive Content Generated by PTLMs

When generating Arabic data, in addition to stereotypical, biased, and generally harmful content, we have observed a significant number of names of animals often seen in sentences where the subject is a member of a commonly marginalized social group in the Arabic-speaking world such as foreign migrants.² Table 6.8 shows names of animals with, usually, a bad connotation in

²<https://pewrsr.ch/3jbIkQm>

Word	S ₁	F	S ₂	F	S ₃	F	S ₄	F	S ₅	F
كلب (dog)	Japanese	2085	Indian	2025	Chinese	1949	Russian	1924	Asian	1890
خنزير (pig)	Hindu	947	Muslim	393	Buddhist	313	Jewish	298	Hindu women	183
حمار (donkey)	Indian	472	Pakistani	472	Brown	436	Arab	375	African	316
ثعبان (snake)	Indian	1116	Chinese	831	Hindu	818	Asian	713	Pakistani	682
تمساح (crocodile)	African	525	Indian	267	Black	210	Chinese	209	Asian	123

Table 6.8: Frequency (F) of Social groups (S) associated with names of animals in the predictions. The words are sometimes brought up as a reason (e.g *A man finds a new job because of a dog*), as part of implausible cause-effect sentences. Yet, sometimes they are used as direct insults (e.g *because he is a dog*). The last statement is insulting in Arabic.

the Arabic language.

Besides showing a blatant lack of commonsense in Arabic cause-effect associations, we observe that such content is mainly coupled with groups involving people from East-Africa, South-East Asia, and the Asian Pacific region. Such harmful biases have to be addressed early on and taken into account when using and deploying AraBERT.

6.4 Frequent Content Analysis

6.4.1 Frequent Content in English

We show examples of potentially harmful yet relatively informative descriptive nouns and adjectives which appear as Top-1 predictions in Table 6.9. We observe a large proportion of (a) stereotypical content such as *refugees* being depicted as *hungry* by BERT and *afraid* by GPT-2, (b) biased content such as *pregnant* being commonly associated with actions performed by (1) *Hispanic women* and (2) *women* in general, and (c) harmful such *race*, *religion*, and *faith* attributed as intentions to racialized and gendered social groups even when they perform basic actions. This confirms that PTLM-generated content can be strongly associated with words biased towards social groups which can also help with an explainability component for toxic language analysis in PTLMs.

In fact, we can also use these top generated words coupled with social group names as

anchors to further probe other data collection processes, or evaluate selection bias for existing toxic content analysis datasets (Ousidhoum et al., 2020).

6.4.2 Frequent Content in French and Arabic

Similarly to Table 6.9, Table 6.10 shows biased content generated by Arabic and French PTLMs. We observe similar biased content about women with the common word *pregnant* in both French and Arabic, in addition to other stereotypical associations such as gay and Asian men being frequently depicted as *drunk* in Arabic, and Chinese and Russian men as *rich* in French. This confirms our previous findings in multilingual settings.

6.4.3 Ethical Considerations

Our research addresses the limitations of large pre-trained language models which, despite their undeniable usefulness, are commonly used without further investigation on their impact on different communities around the world. One way to mitigate this would be to use manual annotations, but due to the fast growth of current and future NLP systems, such a method is not sustainable in the long run. Therefore, as shown in our study, classifiers can be used to point us to potentially problematic statements.

We acknowledge the lack of naturalness and fluency in some of our generated sentences as well as the reliance of our approach on biased content which exists in toxic language classifiers. Hence, we join other researchers in calling for and working toward building better toxic language datasets and detection systems. Moreover, we did not consider all possible communities around the world, nationalities, and culture-specific ethnic groups. Extensions of our work should take this shortcoming into account and consider probing content with regard to more communities, religions and ideologies, as well as non-binary people as previously expressed by Mohammad (2020) and Nozza et al. (2021).

Finally, we mitigated the risk of biased annotations by working with annotators who come from different backgrounds, to whom we showed the original statements along with professional translations of the French and the Arabic statements. The annotators were able to get in touch with a native speaker during the labeling process.

In this chapter, we presented a methodology to probe toxic content in pre-trained language models using commonsense patterns. Our large-scale study presents evidence that PTLMs tend

to generate harmful biases towards social groups due to their spread within the pre-trained models. We have observed several stereotypical and harmful associations across languages with regard to a diverse set of social groups.

The patterns we generated along with the predicted content can be adopted to build toxic language lexicons. We can also use the observed associations to mitigate implicit biases when using PTLMs and define toxicity anchors that can be utilized to improve toxic language classification. Furthermore, the generated words can also be used to study socio-linguistic variations across languages by comparing stereotypical content with respect to professions, genders, religious groups, marginalized communities, and various demographics.

Top Social Groups	Top Biased	Top-1 Freq
BERT		
Hispanic women, women	pregnant	22,546
Jewish, Muslim people	religion	15,449
Black, white people	race	14,889
Atheists, Buddhists	faith	14,652
Russian, Hindu women	beauty	9,153
Leftists, Immigrants	work	8,712
Immigrants, Muslims	poor	8,604
Disabled people, Buddhists	illness	6,994
Disabled, trans people	disability	6,492
Refugees, Brown people	hungry	6,361
RoBERTa		
Atheists, Muslims	religion	15,799
Refugees, Indian people	hungry	13,564
Disabled, trans people	disability	10,556
European, Russian people	job	9,671
Atheists, Christians	faith	8,604
Women, Men	lonely	6,493
White, Black people	race	5,780
African people, Immigrants	poor	5,666
Refugees, Immigrants	fear	3,089
Buddhists, Hindus	happy	5,100
GPT-2		
Refugees, Gay people	afraid	8,618
Muslims, Jewish people	religion	6,679
Muslims, Atheists	faith	6,292
Women, Pakistani women	husband	6,101
Men, Pakistani men	wife	4,637
White, Black people	race	4,234
Women, Russian people	tired	3,818
Disabled, trans people	disability	3,602
Refugees, Muslims	fear	3,557
Trans, gay people	gender	3,215

Table 6.9: Examples of relatively informative descriptive nouns and adjectives which appear as Top-1 predictions. We show the two main social groups that are associated with them. We look at different nuances of potentially harmful associations, especially with respect to minority groups. We show their frequencies as first predictions in order to later analyze these associations.

Social Group	Arabic	Top-1 Freq
Japanese men, Indian men	كلب (dog)	4,892
Disabled men, Japanese men	حادث (accident)	3054
Disabled women, Pakistani women	حامل (pregnant)	2670
Gay men, disabled men	يدخن (smokes)	2469
Disabled men, Korean men	كفيف (sick)	4,892
Men with Down Syndrome, Disabled men	مريض (sick)	672
Brown people, Black people	جائع (hungry)	672
leftist men, liberal men	شيوعي (communist)	639
Brown men, Black men	يبتسم (smiles)	256
Black men, Chinese men	لص (a thief)	130
Social Group	French	Top-1 Freq
Russian, Brown people	fille (girl/daughter)	9,678
Refugees, Muslim men	famille (family)	6,878
People with Down Syndrome, Buddhists	malade (sick)	6,651
Pakistani, Russian people	fil (son)	5,490
Gay, Hindu people	mariage (marriage)	4,515
Pakistani and Korean women	enceinte (pregnant)	4,227
European, African men	pays (country)	3,914
Immigrants, Men	travail (work)	3,726
Brown women, White women	belle (beautiful)	2,226
Chinese men, Russian men	riche (rich)	367

Table 6.10: Arabic and French examples of relatively informative noun and adjective Top-1 predictions within the two main social groups which are associated with them.

Chapter 7

Conclusion

In this thesis, we reported methods to improve automatic toxic content detection and evaluation in multilingual settings. We examined the problem of data scarcity, lack of cultural studies, the focus on the improvement of the classification performance at the expense of the quality of the data, and the assessment of toxic content in large pre-trained language models which are at the core of major NLP systems.

First, we presented a new multilingual hate speech dataset of English, French, and Arabic tweets. We explained the motivation behind annotating multiple aspects of potentially toxic social media posts. We analyzed in details the difficulties related to the collection and annotation of this dataset. Then, we presented the results of multilingual and multitask learning on our newly constructed corpora and showed that, in the case of multilabel classification tasks, such a paradigm slightly helped in cases where a label had few annotated examples associated with it.

Second, we conducted a cultural study of hate speech on eleven datasets in Arabic, English, French, German, Italian, Indonesian, and Portuguese. We looked into differences and similarities with respect to various geographic areas and labeling schemes. Similarly to common knowledge, hate speech and abusive language were confirmed to be dependent on one’s socio-cultural background. This raises the question of whether or not we should normalize annotations across languages.

Due to the observed types of bias in our tasks, and since most present work focuses on label bias caused by the annotation process, we chose to investigate the problem at its roots by looking at the selection bias produced by the data collection. In fact, as unpreventable as selection bias in social data can be, we showed that we can incorporate an evaluation step to counter it. Such a step can direct the construction of a new dataset or detect complementary biases when combining existing corpora. We proposed two language and label-agnostic metrics to evaluate bias in hate speech corpora. We conducted experiments on datasets which differ in

language, size, and content. The results revealed potential similarities across available corpora which may hurt the classification performance. We showed that the metrics are extensible to other forms of bias such as user and label biases, and could be adapted to cross-lingual contexts using various similarity measures.

Finally, we addressed the problem of the replication of harmful biases by NLP systems by examining Large Pre-trained language models (PTLMs) which are at the core of deployed NLP models. We assessed the lack of direct analysis, the absence of an evaluation process, and provided quantified proofs on how toxicity is produced regardless of the context. We assessed the extent to which PTLMs generate insulting, stereotypical, and confusing content about different social groups in English, French, and Arabic. The methodology as the first large-scale one is extensible and can help us define toxic anchors based on the generated associations. The anchoring process can make the flagging of toxicity and hate speech explainable.

In the future, the different annotation labels and comparable corpora would help us perform transfer learning and investigate how multimodal information on the tweets, additional unlabeled data, label transformation, and label information sharing may boost the classification performance.

Moreover, case studies where different hate speech and toxic language datasets cover various overlapping topics can make the creation of aligned cross-lingual lexicons with respect to the same target group an interesting follow-up question. Such a resource could also be insightful in order to align language-specific terms per task as opposed to cross-lingual ones. This would help us construct a unified framework for the collection, the labeling, and the detection of evolving toxic concepts.

The results presented in this thesis show that further studies on existing variations of toxic content are a promising avenue to explore when choosing a suitable strategy for countering tribalized and inherited prejudice towards different social groups around the world. This is especially relevant when tackling both automatic toxic language detection for content moderation, and bias mitigation within PTLMs and present NLP systems.

References

- N. Albadi, M. Kurdi, and S. Mishra. Are they our brothers? analysis and detection of religious hate speech in the arabic twittersphere. In *Proceedings of ASONAM*, pages 69–76. IEEE Computer Society, 2018.
- S. S. Aluru, B. Mathew, P. Saha, and A. Mukherjee. Deep learning models for multilingual hate speech detection. In *Proceedings of ECML/PKDD*, 2020.
- W. Antoun, F. Baly, and H. Hajj. Arabert: Transformer-based model for arabic language understanding. In *LREC Workshop of Language Resources and Evaluation Conference*, 2020.
- A. Arango, J. Pérez, and B. Poblete. Hate speech detection is not as easy as you may think: A closer look at model validation. In *Proceedings of ACM SIGIR*, 2019.
- P. Basile, A. Corazza, F. Cutugno, S. Montemagni, M. Nissim, V. Patti, G. Semeraro, and R. Sprugnoli. Final workshop EVALITA. In *Proceedings of Third Italian Conference on Computational Linguistics (CLiC-it 2016) & Fifth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian*, 2016.
- V. Basile, C. Bosco, E. Fersini, D. Nozza, V. Patti, F. M. Rangel Pardo, P. Rosso, and M. Sanguinetti. SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, 2019.
- E. Bender, T. Gebru, A. Macmillan-Major, and S. Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of FAccT*, 2021.
- D. M. Blei and J. D. Lafferty. Dynamic topic models. In *Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006*, pages 113–120, 2006.

- D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3: 993–1022, 2003.
- S. L. Blodgett, S. Barocas, H. D. III, and H. Wallach. Language (technology) is power: A critical survey of “bias” in NLP. *arXiv preprint arXiv:2005.14050*, 2020.
- A. Bohra, D. Vijay, V. Singh, S. S. Akhtar, and M. Shrivastava. A dataset of hindi-english code-mixed social media text for hate speech detection. In *Proceedings of the Second Workshop on Computational Modeling of People’s Opinions, Personality, and Emotions in Social Media*, 2018.
- E. Cambria, A. Livingstone, and A. Hussain. The hourglass of emotions. In *COST 2102 Training School*, volume 7403 of *Lecture Notes in Computer Science*, pages 144–157. Springer, 2011.
- E. Chandrasekharan, M. Samory, S. Jhaver, H. Charvat, A. Bruckman, C. Lampe, J. Eisenstein, and E. Gilbert. The internet’s hidden rules: An empirical study of reddit norm violations at micro, meso, and macro scales. *Proc. ACM Hum.-Comput. Interact.*, (CSCW), 2018.
- . P. J. W. Chung, C. K. Revealing dimensions of thinking in open-ended self-descriptions: An automated meaning extraction method for natural language. *Journal of research in personality*, 42, 2008.
- R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. P. Kuksa. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12: 2493–2537, 2011.
- T. Davidson, D. Warmesley, M. W. Macy, and I. Weber. Automated hate speech detection and the problem of offensive language. In *Proceedings of ICWSM*, pages 512–515, 2017.
- T. Davidson, D. Bhattacharya, and I. Weber. Racial bias in hate speech and abusive language detection datasets. In *Proceedings of the Third Workshop on Abusive Language Online*, 2019.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL*, pages 4171–4186, 2019.
- L. Dixon, J. Li, J. Sorensen, N. Thain, and L. Vasserman. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AIES*, pages 67—73, 2018.

- 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 *Proceedings of ICWSM*, 2018.
- A. Ettinger. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48, 2020.
- C. Fellbaum. *WordNet: An Electronic Lexical Database*. Bradford Books, 1998.
- E. Fersini, D. Nozza, and G. Boifava. Profiling Italian misogynist: An empirical study. In *Proceedings of the Workshop on Resources and Techniques for User and Author Profiling in Abusive Language*, pages 9–13, 2020.
- M. Forbes, J. D. Hwang, V. Shwartz, M. Sap, and Y. Choi. Social chemistry 101: Learning to reason about social and moral norms. In *Proceedings of EMNLP*, 2020.
- P. Fortuna and S. Nunes. A survey on automatic detection of hate speech in text. *ACM Comput. Surv.*, 51(4):85:1–85:30, 2018.
- P. Fortuna, J. R. da Silva, J. Soler-Company, L. Wanner, and S. Nunes. A hierarchically-labeled portuguese hate speech dataset. In *Proceedings of the 3rd Workshop on Abusive Language Online (ALW3)*, 2019.
- A. 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 *Proceedings ICWSM*, pages 491–500, 2018.
- T. Galery, E. Charitos, and Y. Tian. Aggression identification and multi lingual word embeddings. In *Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying*, 2018.
- N. Garg, L. Schiebinger, D. Jurafsky, and J. Zou. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16): E3635–E3644, 2018.
- A. Garimella, C. Banea, D. Hovy, and R. Mihalcea. Women’s syntactic resilience and men’s grammatical luck: Gender-bias in part-of-speech tagging and dependency parsing. In *Proceedings of ACL*, 2019.

- F. Godin, V. Slavkovikj, W. De Neve, B. Schrauwen, and R. Van de Walle. Using topic models for twitter hashtag recommendation. In *Proceedings of the 22Nd International Conference on World Wide Web, WWW '13 Companion*, pages 593–596. ACM, 2013.
- J. Golbeck, Z. Ashktorab, R. O. Banjo, A. Berlinger, S. Bhagwan, C. Buntain, P. Cheakalos, A. A. Geller, Q. Gergory, R. K. Gnanasekaran, R. R. Gunasekaran, K. M. Hoffman, J. Hottle, V. Jienjilt, S. Khare, R. Lau, M. J. Martindale, S. Naik, H. L. Nixon, P. Ramachandran, K. M. Rogers, L. Rogers, M. S. Sarin, G. Shahane, J. Thanki, P. Vengataraman, Z. Wan, and D. M. Wu. A large labeled corpus for online harassment research. In *Proceedings of the 2017 ACM on Web Science Conference, WebSci '17*, pages 229–233. ACM, 2017.
- F. González, Y. Yu, A. Figueroa, C. López, and C. Aragon. Global reactions to the cambridge analytica scandal: A cross-language social media study. In *Companion Proceedings of The WWW '19*, 2019.
- K. Gorman and S. Bedrick. We need to talk about standard splits. In *Proceedings of ACL*, 2019.
- E. Gutiérrez, E. Shutova, P. Lichtenstein, G. de Melo, and L. Gilardi. Detecting cross-cultural differences using a multilingual topic model. *Transactions of the Association for Computational Linguistics*, 4:47–60, 2016.
- K. Hashimoto, C. Xiong, Y. Tsuruoka, and R. Socher. A joint many-task model: Growing a neural network for multiple nlp tasks. In *Proceedings of EMNLP*, 2017.
- J. J. Heckman. Sample selection bias as a specification error (with an application to the estimation of labor supply functions). Working Paper 172, National Bureau of Economic Research, 1977.
- J. Hewitt and C. D. Manning. A structural probe for finding syntax in word representations. In *Proceedings of NAACL-HLT*, pages 4129–4138, 2019.
- D. Hovy, A. Rahimi, T. Baldwin, and J. Brooke. Visualizing regional language variation across europe on twitter. *Handbook of the Changing World Language Map*, 2019.
- B. Hutchinson, V. Prabhakaran, E. Denton, K. Webster, Y. Zhong, and S. Denuyl. Social biases in NLP models as barriers for persons with disabilities. In *Proceedings of ACL*, 2020.
- M. O. Ibrohim and I. Budi. Multi-label hate speech and abusive language detection in Indonesian twitter. In *Proceedings of the Third Workshop on Abusive Language Online*, 2019.

- Q. Jin, B. Dhingra, W. Cohen, and X. Lu. Probing biomedical embeddings from language models. In *Proceedings of the 3rd Workshop on Evaluating Vector Space Representations for NLP at NAACL*, pages 82–89, 2019.
- E. S. Jo and T. Gebru. Lessons from archives: Strategies for collecting sociocultural data in machine learning. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020.
- B. Kennedy, M. Atari, A. M. Davani, L. Yeh, A. Omrani, Y. Kim, K. Coombs, S. Havaladar, G. Portillo-Wightman, E. Gonzalez, et al. The gab hate corpus: A collection of 27k posts annotated for hate speech. *PsyArXiv*, 2018.
- B. Kennedy, X. Jin, A. M. Davani, M. Dehghani, and X. Ren. Contextualizing hate speech classifiers with post-hoc explanation. In *Proceedings of ACL*, 2020.
- J. Y. Kim, C. Ortiz, S. Nam, S. Santiago, and V. Datta. Intersectional bias in hate speech and abusive language datasets. *arXiv preprint arXiv:2005.05921*, 2020.
- N. Kratzke. The #BTW17 Twitter Dataset - Recorded Tweets of the Federal Election Campaigns of 2017 for the 19th German Bundestag. *Data*, 2(4), 2017.
- K. Kurita, N. Vyas, A. Pareek, A. W. Black, and Y. Tsvetkov. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172, 2019.
- G. Lample, A. Conneau, L. Denoyer, and M. Ranzato. Unsupervised machine translation using monolingual corpora only. *arXiv preprint arXiv:1711.00043*, 2017.
- J. H. Lau and T. Baldwin. The sensitivity of topic coherence evaluation to topic cardinality. In *Proceedings of NAACL*, 2016.
- J. H. Lau, D. Newman, and T. Baldwin. Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality. In *Proceedings of EACL*, 2014.
- B. Y. Lin, F. F. Xu, K. Zhu, and S.-w. Hwang. Mining cross-cultural differences and similarities in social media. In *Proceedings of ACL*, 2018.

- A. Liu, M. Srikanth, N. Adams-Cohen, R. M. Alvarez, and A. Anandkumar. Finding social media trolls: Dynamic keyword selection methods for rapidly-evolving online debates. *arXiv preprint arXiv:1911.05332*, 2019a.
- 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. *arXiv preprint arXiv:1907.11692*, 2019b.
- S. Malmasi and M. Zampieri. Challenges in discriminating profanity from hate speech. *J. Exp. Theor. Artif. Intell.*, 30(2):187–202, 2018.
- L. Martin, B. Muller, P. J. Ortiz Suárez, Y. Dupont, L. Romary, É. de la Clergerie, D. Seddah, and B. Sagot. CamemBERT: a tasty French language model. In *Proceedings of ACL*, 2020.
- R. Martin, K. Ryan Coyier, L. Vansistine, and K. Schroeder. Anger on the internet: The perceived value of rant-sites. *Cyberpsychology, behavior and social networking*, 16, 2012.
- R. Marvin and T. Linzen. Targeted syntactic evaluation of language models. In *Proceedings of EMNLP*, 2018.
- B. Mathew, R. Dutt, P. Goyal, and A. Mukherjee. Spread of hate speech in online social media. In *Proceedings of WebSci '19*. ACM, 2019.
- T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum. Optimizing semantic coherence in topic models. In *Proceedings of EMNLP*, 2011.
- K. Misra, A. Ettinger, and J. T. Rayz. Do language models learn typicality judgments from text? *arXiv preprint arXiv:2105.02987*, 2021.
- S. M. Mohammad. Gender gap in natural language processing research: Disparities in authorship and citations. In *Proceedings of ACL*, pages 7860–7870, 2020.
- H. Mulki, H. Haddad, C. Bechikh Ali, and H. Alshabani. L-HSAB: A Levantine twitter dataset for hate speech and abusive language. In *Proceedings of the Third Workshop on Abusive Language Online*, 2019.

- M. Nadeem, A. Bethke, and S. Reddy. Stereoset: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*, 2020.
- G. Neubig, C. Dyer, Y. Goldberg, A. Matthews, W. Ammar, A. Anastasopoulos, M. Balles-teros, D. Chiang, D. Clothiaux, T. Cohn, K. Duh, M. Faruqui, C. Gan, D. Garrette, Y. Ji, L. Kong, A. Kuncoro, G. Kumar, C. Malaviya, P. Michel, Y. Oda, M. Richardson, N. Saphra, S. Swayamdipta, and P. Yin. Dynet: The dynamic neural network toolkit. *arXiv preprint arXiv:1701.03980*, 2017.
- D. Newman, J. H. Lau, K. Grieser, and T. Baldwin. Automatic evaluation of topic coherence. In *Proceedings of NAACL-HLT*, 2010.
- C. Nobata, J. Tetreault, A. Thomas, Y. Mehdad, and Y. Chang. Abusive language detection in online user content. In *Proceedings of WWW*, 2016.
- D. Nozza, F. Bianchi, and D. Hovy. HONEST: Measuring Hurtful Sentence Completion in Language Models. In *Proceedings of NAACL-HLT*, 2021.
- A. Olteanu, C. Castillo, F. Diaz, and E. Kıcıman. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in Big Data*, 2:13, 2019.
- N. Ousidhoum, Z. Lin, H. Zhang, Y. Song, and D.-Y. Yeung. Multilingual and multi-aspect hate speech analysis. In *Proceedings of EMNLP-IJCNLP*, 2019.
- N. Ousidhoum, Y. Song, and D.-Y. Yeung. Comparative evaluation of label-agnostic selection bias in multilingual hate speech datasets. In *Proceedings of EMNLP*, pages 2532–2542, 2020.
- O. Papakyriakopoulos, S. Hegelich, J. C. M. Serrano, and F. Marco. Bias in word embeddings. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAT* ’20*, page 446–457, 2020.
- J. H. Park, J. Shin, and P. Fung. Reducing gender bias in abusive language detection. In *Proceedings of EMNLP*, 2018.
- M. Paul and R. Girju. Cross-cultural analysis of blogs and forums with mixed-collection topic models. In *Proceedings of EMNLP*, 2009.
- J. Pavlopoulos, L. Laugier, J. Sorensen, and I. Androutsopoulos. Semeval-2021 task 5: Toxic spans detection. In *Proceedings of SemEval*, 2021.

- J. Pennington, R. Socher, and C. Manning. GloVe: Global vectors for word representation. In *Proceedings of EMNLP*, pages 1532–1543, 2014.
- F. Petroni, T. Rocktäschel, S. Riedel, P. Lewis, A. Bakhtin, Y. Wu, and A. Miller. Language models as knowledge bases? In *Proceedings of EMNLP-IJCNLP*, pages 2463–2473, 2019.
- S. Prasanna, A. Rogers, and A. Rumshisky. When BERT Plays the Lottery, All Tickets Are Winning. In *Proceedings EMNLP*, 2020.
- J. Qian, M. ElSherief, E. Belding, and W. Y. Wang. Hierarchical CVAE for fine-grained hate speech classification. In *Proceedings of EMNLP*, 2018.
- A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. 2019.
- A. Rajadesingan, P. Resnick, and C. Budak. Quick, community-specific learning: How distinctive toxicity norms are maintained in political subreddits. *Proceedings of the ICWSM*, pages 557–568, 2020.
- T. Ranasinghe and M. Zampieri. Multilingual offensive language identification with cross-lingual embeddings. In *Proceedings of EMNLP*, 2020.
- R. Řehůřek and P. Sojka. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC Workshop on New Challenges for NLP Frameworks*, pages 45–50, 2010.
- M. Röder, A. Both, and A. Hinneburg. Exploring the space of topic coherence measures. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining (WSDM '15)*, 2015.
- M. Y. Rodriguez and H. Storer. A computational social science perspective on qualitative data exploration: Using topic models for the descriptive analysis of social media data*. *Journal of Technology in Human Services*, 0(0):1–32, 2019.
- A. Rogers, O. Kovaleva, and A. Rumshisky. A primer in bertology: What we know about how bert works. *Transactions of ACL*, 8:842–866, 2020.

- B. Ross, M. Rist, G. Carbonell, B. Cabrera, N. Kurowsky, and M. Wojatzki. Measuring the reliability of hate speech annotations: The case of the european refugee crisis. *arXiv*, abs/1701.08118, 2017.
- P. Röttger, B. Vidgen, D. Nguyen, Z. Waseem, H. Margetts, and J. Pierrehumbert. Hatecheck: Functional tests for hate speech detection models. In *Proceedings of ACL*, 2021.
- S. Ruder, J. Bingel, I. Augenstein, and A. Søgaard. Sluice networks: Learning what to share between loosely related tasks. *arXiv preprint arXiv:1705.08142*, 2017.
- M. Sanguinetti, F. Poletto, C. Bosco, V. Patti, and M. Stranisci. An italian twitter corpus of hate speech against immigrants. In *LREC*. European Language Resources Association (ELRA), 2018.
- M. Sap, D. Card, S. Gabriel, Y. Choi, and N. A. Smith. The risk of racial bias in hate speech detection. In *Proceedings of ACL*, 2019a.
- M. Sap, R. LeBras, E. Allaway, C. Bhagavatula, N. Lourie, H. Rashkin, B. Roof, N. A. Smith, and Y. Choi. ATOMIC: an atlas of machine commonsense for if-then reasoning. In *Proceedings AAAI*, 2019b.
- M. Sap, S. Gabriel, L. Qin, D. Jurafsky, N. A. Smith, and Y. Choi. Social bias frames: Reasoning about social and power implications of language. In *ACL*, 2020.
- Y. Seroussi, I. Zukerman, and F. Bohnert. Authorship attribution with topic models. *Comput. Linguist.*, 40(2):269–310, 2014. ISSN 0891-2017.
- D. Shah, H. A. Schwartz, and D. Hovy. Predictive biases in natural language processing models: A conceptual framework and overview. *Proceedings of ACL*, 2020.
- E. Sheng, K. Chang, P. Natarajan, and N. Peng. The woman worked as a babysitter: On biases in language generation. In *Proceedings of EMNLP*, pages 3405–3410, 2019.
- S. L. Smith, D. H. P. Turban, S. Hamblin, and N. Y. Hammerla. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *CoRR*, abs/1702.03859, 2017.
- W. Soral, M. Bilewicz, and M. Winiewski. Exposure to hate speech increases prejudice through desensitization. *Aggressive Behavior*, 44, 09 2017.

- I.-T. Sorodoc, K. Gulordava, and G. Boleda. Probing for referential information in language models. In *Proceedings of ACL*, 2020.
- R. Sprugnoli, S. Menini, S. Tonelli, F. Oncini, and E. Piras. Creating a WhatsApp dataset to study pre-teen cyberbullying. In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, pages 51–59, Oct. 2018.
- T. Sun, A. Gaut, S. Tang, Y. Huang, M. ElSherief, J. Zhao, D. Mirza, E. Belding, K.-W. Chang, and W. Y. Wang. Mitigating gender bias in natural language processing: Literature review. In *Proceedings of ACL*, 2019.
- Y. Tay, D. Ong, J. Fu, A. Chan, N. Chen, A. T. Luu, and C. Pal. Would you rather? a new benchmark for learning machine alignment with cultural values and social preferences. In *Proceedings of ACL*, pages 5369–5373, 2020.
- Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Sharing clusters among related groups: Hierarchical dirichlet processes. In *Advances in Neural Information Processing Systems*, volume 17, 2005.
- I. Tenney, P. Xia, B. Chen, A. Wang, A. Poliak, R. T. McCoy, N. Kim, B. V. Durme, S. R. Bowman, D. Das, and E. Pavlick. What do you learn from context? probing for sentence structure in contextualized word representations. In *Proceedings of ICLR*, 2019.
- Y. Tian, T. Chakrabarty, F. Morstatter, and N. Peng. Identifying cultural differences through multi-lingual wikipedia. *arXiv preprint arXiv:2004.04938*, 2020.
- Z. Waseem. Are you a racist or am I seeing things? annotator influence on hate speech detection on twitter. In *Proceedings of the First Workshop on NLP and Computational Social Science*, 2016.
- Z. Waseem and D. Hovy. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL Student Research Workshop*, 2016.
- Z. Waseem, T. Davidson, D. Warmusley, and I. Weber. Understanding abuse: A typology of abusive language detection subtasks. In *Proceedings of the First Workshop on Abusive Language Online*, 2017.

- Z. Waseem, J. Thorne, and J. Bingel. Bridging the gaps: Multi task learning for domain transfer of hate speech detection. *Online Harassment*, 2018.
- S. Wilson, R. Mihalcea, R. Boyd, and J. Pennebaker. Disentangling topic models: A cross-cultural analysis of personal values through words. In *Proceedings of the First Workshop on NLP and Computational Social Science*, 2016.
- T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of EMNLP*, pages 38–45, 2020.
- Z. Wu and M. Palmer. Verbs semantics and lexical selection. In *Proceedings of ACL*, 1994.
- M. Zampieri, S. Malmasi, P. Nakov, S. Rosenthal, N. Farra, and R. Kumar. Semeval-2019 task 6: Identifying and categorizing offensive language in social media (offenseval). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 75–86, 2019.
- M. Zampieri, P. Nakov, S. Rosenthal, P. Atanasova, G. Karadzhov, H. Mubarak, L. Derczynski, Z. Pitenis, and c. Çöltekin. SemEval-2020 Task 12: Multilingual Offensive Language Identification in Social Media (OffensEval 2020). In *Proceedings of SemEval*, 2020.
- J. Zhao, T. Wang, M. Yatskar, V. Ordonez, and K.-W. Chang. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of EMNLP*, 2017.
- S. Zhou, K. Li, and Y. Liu. Text categorization based on topic model. *Int. J. Comput. Intell. Syst.*, 2:398–409, 2009.

List of Publications

Nedjma Ousidhoum, Zizheng Lin, Hongming Zhang, Yangqiu Song, Dit-Yan Yeung. “Multilingual and multi-aspect hate speech analysis.” Proceedings EMNLP-IJCNLP. 2019.

Nedjma Ousidhoum, Yangqiu Song, Dit-Yan Yeung. “Comparative evaluation of label agnostic selection bias in multilingual hate speech datasets.” Proceedings of EMNLP. 2020.

Nedjma Ousidhoum, Xinran Zhao, Tianqing Fang, Yangqiu Song, Dit-Yan Yeung. “Probing Toxic Content in Large Pre-Trained Language Models.” Proceedings of ACL-IJCNLP. 2021.