

# Data Bootstrapping Approaches to Improve Low Resource Abusive Language Detection for Indic Languages

Mithun Das  
mithundas@iitkgp.ac.in  
Indian Institute of Technology  
Kharagpur  
West Bengal, India

Somnath Banerjee  
som.iitkgpcse@kgpian.iitkgp.ac.in  
Indian Institute of Technology  
Kharagpur  
West Bengal, India

Animesh Mukherjee  
animeshm@cse.iitkgp.ac.in  
Indian Institute of Technology  
Kharagpur  
West Bengal, India

## ABSTRACT

Abusive language is a growing concern in many social media platforms. Repeated exposure to abusive speech has created physiological effects on the target users. Thus, the problem of abusive language should be addressed in all forms for online peace and safety. While extensive research exists in abusive speech detection, most studies focus on English. Recently, many smearing incidents have occurred in India, which provoked diverse forms of abusive speech in online space in various languages based on the geographic location. Therefore it is essential to deal with such malicious content. In this paper, to bridge the gap, we demonstrate a large-scale analysis of multilingual abusive speech in Indic languages. We examine different interlingual transfer mechanisms and observe the performance of various multilingual models for abusive speech detection for eight different Indic languages. We also experiment to show how robust these models are on adversarial attacks. Finally, we conduct an in-depth error analysis by looking into the models' misclassified posts across various settings. We have made our code and models public for other researchers<sup>1</sup>.

## CCS CONCEPTS

- **Computing methodologies** → **Natural language processing**;
- **Social and professional topics** → **Censorship**.

## KEYWORDS

Abusive language, multilingual, detection, social media

## 1 INTRODUCTION

Social media platforms (such as Twitter, Facebook, etc.) have connected billions of people at different levels and allowed them to share ideas among themselves instantly. On the one hand, it has led to the exchange of thoughts and massive expansion of social networks; on the other hand, these platforms have been used to spread propaganda, violence, and abuse against users based on gender, religion, race, geographic location, etc. [11]. Not only such abusive behavior can lead to the traumatization of the victims by affecting them psychologically [37], but these can also ignite social tensions and affect the stature of the platforms which host them [36]. Further, widespread usage of such content can also have implications in the offline world: violent hate crimes, youth suicides, mass shootings, and extremist recruitment [17].

To mitigate the spread of such abusive content, these platforms have come up with specific guidelines<sup>2</sup> that need to be followed

by the users of these platforms. Failure to follow these guidelines may result in the deletion of the post or suspension of the user's account. To reduce the hateful content from these platforms, they employ moderators [28] to review posts manually and keep the forum healthy and friendly. However, this moderation technique is limited by the moderators' speed, diction, ability to understand the evolution of vernacular, and familiarity with multilingual content. In addition, many moderators complain about psychological effects induced due to moderation of such abusive content[1]. Besides, due to the sheer volume and velocity of data streaming, it is an ambitious attempt to filter all the posts manually and screen out such hostile content. Thus, automatic detection of such abusive content is incredibly crucial and unavoidable.

It has already been eyed that Facebook actively removed a large portion of malicious content from their platform even before users flagged it<sup>3</sup>. Nevertheless, the limitation is that these platforms can identify such detrimental content in certain major languages such as English, Spanish, etc. To that end, several studies have been performed for automated detection of abusive content, concentrating predominantly on the English language. Hence, an effort is needed to determine and mitigate such malicious content in low-resource languages.

In the last few years, there have been a series of incidents in India, such as smearing movements against famous political leaders<sup>4</sup>, celebrities<sup>5</sup>, and social media personalities, online anti-religious propaganda<sup>6</sup>, cyber harassment<sup>7</sup>, etc. So, to deal with such malicious content, automated systems are much required to keep the online ecosystem healthy. India has more than 1.3 billion people, having the highest number of users on Facebook<sup>8</sup>, YouTube<sup>9</sup> and the third-highest number of users on Twitter<sup>10</sup>. Besides, the country has 22 recognized languages, which are spoken in various parts of it<sup>11</sup>. Due to the most extensive language diversity and their usage on social media, detecting such abusive content in all languages becomes challenging. There are primarily two reasons for this –

<sup>3</sup><https://time.com/5739688/facebook-hate-speech-languages/>

<sup>4</sup><https://nenow.in/top-news/resign-pm-modi-trends-on-twitter-as-india-sees-worst-covid-19-wave.html>

<sup>5</sup><https://www.koimoi.com/bollywood-news/sushant-singh-rajput-news-rhea-chakraborty-was-harassed-by-various-agencies-satish-maneshinde-demands-cbi-findings>

<sup>6</sup><https://www.news18.com/news/buzz/indians-wants-to-boycott-myntna-for-old-anti-hindu-poster-it-didnt-even-make-4115855.html>

<sup>7</sup><https://timesofindia.indiatimes.com/city/mumbai/cyber-harassment-cases-see-upswing-in-pandemic/articleshow/88842765.cms>

<sup>8</sup><https://worldpopulationreview.com/country-rankings/facebook-users-by-country>

<sup>9</sup><https://www.globalmediainsight.com/blog/youtube-users-statistics/>

<sup>10</sup><https://www.statista.com/statistics/242606/number-of-active-twitter-users-in-selected-countries/>

<sup>11</sup><https://www.universal-translation-services.com/recognized-languages-in-india/>

<sup>1</sup><https://github.com/hate-alert/IndicAbusive>

<sup>2</sup><https://help.twitter.com/en/rules-and-policies/hateful-conduct-policy>

annotators need to have diverse language expertise and should have knowledge in abusive content analysis.

Further, a recent trend on social media platforms is that people write non-English languages using English characters and switch among two or more languages in the same conversation. This phenomenon is called code-mixing (or code-switching), where linguistic units such as phrases, words, or morphemes of one language are embedded in an utterance of another language. Code-mixing allows ease of communication among speakers by furnishing a more comprehensive set of phrases and expressions. However, this has also made the task of developing NLP tools more complex, as emphasized by Chittaranjan et al. [9]. Therefore, there is a need to develop efficient models to detect abusive content in Indic languages in various settings.

This paper performs a large-scale analysis of multilingual abusive speech by investigating multilingual models’ performance in eight different Indic languages. We address the question of data scarcity and language similarity/typology, two under-explored issues in abusive language identification. Inspired by Nag et al. [27], we explore and catalog a variety of strategies for transferring abusive language detection ability across languages, especially in the resource-rich to resource-poor direction, starting from the “each language for itself” (ELFI) criteria. To investigate the degree to which various transfer modes can compensate for gold training instances, we explore multiple scenarios ranging from zero-shot learning, few-shot learning, instance transfer, cross-lingual learning, etc. (more details in section 4). In summary, we observe that

- In ELFI style training, although **m-BERT** is pre-trained in more than 100 languages, **MuRIL** outperforms m-BERT in 7 out of 10 language types as it is pre-trained explicitly in Indian languages.
- The zero-shot model transfer can be beneficial for abusive language detection when the source and the target languages belong to the same language family. For few-shot settings, **AllBOne** is the most effective.
- The model transfer brings better performance than instance transfer (i.e., all instances of a resource-poor language are first translated to a resource-rich language and the training and predictions are done in the resource-rich language) due to the added translation error in instance transfer which reduces the overall toxicity score of an abusive post.
- For low-resource languages, though synthetic silver instances are helpful in detecting abusive content up to some degree, further fine-tuning with gold target instances gives steady improvements.
- Our in-depth error analysis reveals that when a post’s contextual information is limited, implicit, or discriminatory features are present, the model fails to detect such an abusive post.

## 2 RELATED WORKS

### 2.1 Dynamics of online abuse

Online hostility is a context-dependent notion intended to express hatred and threaten an individual or group based on discriminatory

views. Despite the argument that hateful statements ought to be tolerated due to free speech acts, the public expression of hate speech propels the reduction of minority members, and such frequent and repetitive exposure to abusive speech could increase an individual’s outgroup prejudice [35]. Real-world violent events could also lead to increased hatred in online space and vice versa [29]. With the rise of online hate, the research community has a massive responsibility to develop solutions to mitigate online hostility.

### 2.2 Research on abusive speech

The concern of abusive speech has long been studied in the research community. Earlier work on abusive speech attempted to detect abusive users by using lexical, syntactic features extracted from their posts [8]. Over the past few years, research around automated hate speech detection has matured tremendously. Most of the current study consists of diverse but related works. In 2016, Zeerak Waseem and Dirk Hovy [38] contributed a dataset in which thousands of tweets were labeled with racism and sexism markers, and Davidson et al. [13] focused on distinguishing offensive from hate content on Twitter. Using this dataset, the authors examined multiple linguistic features such as character and word n-grams, POS tags, emotion lexicon, and tf-idf vectors with several classifiers such as LR, SVM, decision tree, etc. Although multiple datasets were being published, a major problem was the lack of correlation and reusable datasets across hate speech detection tasks [21]. To address this issue, Founta et al. [15] studied these inconsistencies and drew upon a robust labeling mechanism that attempts to circumvent the overlap among various forms of abusive speech.

With the advent of large datasets, most academic research has moved to data-hungry complex models to improve classifier performance, including deep learning [4] and graph embedding techniques [12]. Pitsilis et al. [30], used deep learning models such as LSTMs to identify the abusive tweets in English and noticed that it was pretty effective in this task. Zhang et al. [39] fused convolutional and gated recurrent networks to enhance the classification performance and had remarkable success on 6 out of 7 datasets used. Recently, transformer based language models such as BERT are becoming immensely popular in several downstream tasks and have outperformed several deep learning models such as CNN-GRU, LSTM, etc., for detecting abusive language [5, 10].

### 2.3 Abusive language detection in Indic languages

In the last few years, several shared tasks, such as Hate-Speech and Offensive Content Identification (HASOC) [25], Dravidian Lang-Tech [7] workshop, TRAC [22], etc., have been organized to develop resources, datasets, and models for abusive speech detection and multiple datasets in Indic languages such as Hindi, Marathi, Tamil, Malayalam, etc. have been made public. The HASOC [25] shared task in Indo-European languages is arguably the most well-known series of competitions. It has been consistently organized from 2019 at the Forum for Information Retrieval (FIRE). The Dravidian Lang-Tech [7] workshop focused on determining the offensive language of the code-mixed dataset in three Dravidian languages, namely, Tamil-English, Malayalam-English, and Kannada-English crawled from social media. In addition, researchers have also developed

several datasets for Bengali [34], code-mixed Hindi [6], Urdu [2, 33], etc., for abusive language detection. However, a limited number of studies have been performed on the effect of zero-shot learning, few-shot learning [31, 32], instance transfer, etc. In our work, we try to fill this critical gap by studying various transfer schemes thus opening up new avenues for future research for abuse detection in Indic languages.

### 3 DATASETS

In this section, we describe the datasets used in this paper. We looked into several datasets for Indic languages for abusive speech detection and attempted to gather all of them. To this purpose, we accumulated datasets in 8 (10 types) different languages from 14 publicly available sources. These datasets differ in their choice of class labels (*offensive*, *hate-speech*, *normal*). Apart from normal posts, we combine all other labels as abusive to pose the problem as a binary classification task. The details about the datasets are presented in Table 1.

**Bengali (Bn):** The dataset shared by Romin et al. [34] is one of the largest in Bengali, consisting of 30K posts, among which 10K posts are hateful, and 20K posts are normal. The dataset has been developed by extracting comments from YouTube and Facebook pages. The author employed 50 annotators and instructed them with proper guidelines to build the dataset. To validate the annotation quality, the author randomly sampled 300 posts for each annotator, manually reannotated them, and found that 91.05% annotation was correct.

**English (En):** The majority of the abusive speech datasets are available in the English language, and out of these, we select only three publicly available popular datasets. The work on automatic hate speech detection by Davidson et al. [13] made public a Twitter dataset consisting of 24k tweets. To curate the dataset, they used a set of lexicons derived from Hatebase.org<sup>12</sup>. Each tweet has been labeled as either hate speech, offensive or normal. Founta et al. [15] shared a large-scale Twitter dataset consisting of more than 100K tweets. Each tweet has been labeled into one of the following categories: hateful, abusive, normal, and spam. The tweets were annotated by 5-20 annotators to maintain the quality of the labels. We ignore the data points annotated as spam from our analysis. The HateXplain dataset introduced by Mathew et al. [26] is a recent addition to the hate speech research community. It contains around 20K posts categorized into three labels: hate speech, offensive, or normal. The dataset is collected from Twitter and Gab.

**Hindi:** For the Hindi language, we found two types of datasets: Devanagari Hindi and code-mixed Hindi. The difference between Devanagari and code-mixed is, though the semantic expression of both types of posts is almost similar, Devanagari Hindi is written using the Devanagari script while code-mixed Hindi is written using English characters. Due to differences in writing style, we treat these languages separately.

- **Devanagari Hindi (Hi):** In 2019, Mandl et al. [24] released a dataset of 5.9K tweets via the HASOC shared task. The posts were labeled into one of the following classes: hate, offensive, profane and normal. The author developed the dataset by crawling tweets and comments from Twitter and

Facebook. Following the previous work, in 2020, Mandl et al. [23] shared another Hindi abusive speech dataset of 3.6K. Recently, the author introduced another dataset [25] of 6.1K due to the huge success of the previously organized shared task. We combine all three datasets to have our final Hindi dataset.

- **Code-mixed Hindi (Hi-En):** Bohra et al. [6] introduced the first code-mixed Hindi hate speech Twitter dataset. Each tweet has been annotated as either hate speech or non-hate speech. According to their annotation guidelines, they considered any kind of abuse as hateful. The final dataset consists of 4.5K tweets, out of which 1.6K tweets are hateful, and the remaining 2.9K are non-hateful.

**Kannada (Ka-En):** Chakravarthi et al. [7] released an offensive language identification dataset in Kannada-English (i.e., code-mixed Kannada). To develop the dataset, the author crawled YouTube comments in 2019. The comments are labeled into one of the following categories: not-offensive, offensive-untargeted, offensive-targeted-individual, offensive-targeted-group, offensive-targeted-other, and not-in-indented-language. The final dataset consists of around 7.7K comments. We removed the data points labeled as not-in-indented-language for our analysis, which was irrelevant. All the not-offensive points are assumed to be in the normal class while all the other points (except not-in-indented-language) are used to form the abusive class.

**Malayalam (Ma-En):** Similar to the Kannada dataset, Chakravarthi et al. [7] made public another code-mixed offensive language detection dataset in Malayalam. The dataset consists of around 20K comments, out of which around 700 comments are offensive. Mandl et al. [23] made public another code-mixed Malayalam dataset consisting 4.9K comments out of which 2.4K comments are abusive.

**Marathi (Mr):** The Marathi dataset shared by Gaikwad et al. [16] consists of 2.4K posts, among which 1.62K posts are labeled as offensive, and the remaining posts are marked as non-offensive. The dataset has been curated by extracting tweets from Twitter by searching common curse words in Marathi and phrases related to politics, entertainment, and sports. The author employed six volunteer annotators who were native speakers of Marathi with ages between 20 and 25 years old and a bachelor's degree.

**Tamil (Ta-En):** Chakravarthi et al. [7] made public another offensive language detection dataset in code-mixed Tamil. The dataset consists of 43K comments, making it one of the most extensive datasets in code-mixed Tamil.

**Urdu:** We also found two types (actual and code-mixed) of datasets for the Urdu language.

- **Actual Urdu (Ur):** For the Urdu language, we found two publicly available datasets. Akhter et al. [2] shared a dataset consisting of 2.1K (1.1K offensive, 1K non-offensive) comments scrapped from YouTube videos. The dataset was manually annotated by three annotators who are native speakers of Urdu. Amjad et al. [3] made public another Urdu abusive dataset crawled from Twitter, consisting of 3.5K tweets, out of which 1.75K tweets are abusive, and the remaining are labeled as non-abusive.
- **Code-mixed Urdu (Ur-En):** For code-mixed Urdu, we use the following two datasets. Khan et al. [18] released a dataset

<sup>12</sup><https://hatebase.org/>

Language	Type	Dataset	Source	Abusive	Normal	Total
Bengali	<i>Bn</i>	<i>Romin et al.</i> [34]	Facebook & YouTube	10,000	20,000	30,000
English	<i>En</i>	<i>Davidson et al.</i> [13]	Twitter	20,620	4,163	24,783
		<i>Founta et al.</i> [15]	Twitter	31,985	53,790	85,775
		<i>HateXplain</i> [26]	Twitter & Gab	11,415	7,814	19,229
Hindi	<i>Hi</i>	<i>Mandl et al.</i> [24]	Twitter & Facebook	3,074	2,909	5,983
		<i>Mandl et al.</i> [23]	Twitter	1,044	2,582	3,626
	<i>Mandl et al.</i> [25]	Twitter	1,938	4,188	6,126	
	<i>Hi-En</i>	<i>Bohra et al.</i> [6]	Twitter	1,661	2,918	4,579
Kannada	<i>Ka-En</i>	<i>Chakravarthi et al.</i> [7]	YouTube	1,465	4,188	5,873
Malayalam	<i>Ma-En</i>	<i>Chakravarthi et al.</i> [7]	YouTube	706	17,697	18,403
		<i>Mandl et al.</i> [23]	YouTube	2,430	2,520	4,950
Marathi	<i>Mr</i>	<i>Gaikwad et al.</i> [16]	Twitter	865	1,611	2,499
Tamil	<i>Ta-En</i>	<i>Chakravarthi et al.</i> [7]	YouTube	12,651	33,684	47,072
Urdu	<i>Ur</i>	<i>Akhter et al.</i> [2]	YouTube	1,109	1,062	2,171
		<i>Amjad et al.</i> [3]	Twitter	1,750	1,750	3,500
	<i>Ur-En</i>	<i>Khan et al.</i> [18]	Twitter	3,575	1,425	5,000
		<i>Rizwan et al.</i> [33]	Twitter	5,349	4,664	10,013
Total	-	-	-	111,637	159,859	271,496

Table 1: Details of the datasets.

of 5K tweets. Each tweet has been annotated into one of the following categories: hate, offensive, neutral. Rizwan et al. [33] made public another code-mixed Urdu dataset having 10K posts containing five different classes: offensive, sexism, religious hate, profane, and normal. Excluding normal, we map all other classes into one label, i.e., abusive.

## 4 METHODOLOGY

### 4.1 Base models

**m-BERT(MB)**: m-BERT [14] is pre-trained on 104 languages with the largest Wikipedia utilizing a masked language modeling (MLM) objective. It is a stack of transformer encoder layers with 12 “attention heads,” i.e., fully connected neural networks augmented with a self-attention mechanism. m-BERT is restricted in the number of tokens it can handle (512 at max). To fine-tune m-BERT, we also add a fully connected layer with the output corresponding to the CLS token in the input. This CLS token output usually holds the representation of the sentence passed to the model. The m-BERT model has been well studied in abusive speech, has already surpassed existing baselines, and stands as a state-of-the-art.

**MuRIL(MU)**: MuRIL [19] stands for Multilingual Representations for Indian Languages and aims to improve interoperability from one language to another. This model uses a BERT base architecture pretrained from scratch utilizing the Wikipedia, Common Crawl, PMINDIA, and Dakshina corpora for 17 Indian languages and their transliterated counterparts.

### 4.2 Interlingual transfer mechanisms

One of the primary interests of transformer-based models is their potential to leverage model transfer via several mechanisms. This can be especially helpful for improving the performance of learning

in low-resource languages like Bengali, Hindi, Urdu, etc. We perform the following tests to evaluate the extent to which language similarity can boost transfer learning performance.

**4.2.1 ELFI (Each language for itself)**. We use data from the same language for training, validation, and testing in this setting. This scenario usually occurs in the real world, where annotated (labeled) datasets are used to create classifications for a specific language. While the labeling costs are potentially high, this provides an idea of the most feasible classification performance.

**4.2.2 Joint training/Cross-lingual training**. In this technique, we combine the datasets of all the languages for training the transformer-based models. The idea is that even though the characters, words used to represent different languages vary, the contextual representation of these abusive posts is the same to good extent. In specific, we consider pre-trained embeddings of all the datapoints from all languages (inclusive of the target language) and test it on the test data of the target language. Thus, it gives an idea of whether jointly training the models can help learn a particular post’s better semantic representation for determining its corresponding label.

**4.2.3 Model transfer**. In this setting, the models are trained with one language (source language) and assessed on another language (target language). In the zero-shot setting, no instances from the target language have been used while training (**MTx0**). In a related few-shot setting, we allow  $n = 32$  and 64 posts per label from the available gold target instances to fine-tune the current models (trained on another language). These are called **MTx32** and **MTx64**. Another extended variant of this model is where for a target language, we use the dataset of all other languages (combined source) to train the models and evaluate their performance on the target language. It would be the case in which we would like to deploy an abusive speech classifier for a target language directly which does

Lang	Model	Accuracy	M-F1
<i>Bn</i>	MB	0.903	0.892
	MU	<b>0.905</b>	<b>0.894</b>
<i>Hi</i>	MB	0.776	0.768
	MU	<b>0.806</b>	<b>0.799</b>
<i>Hi-En</i>	MB	<b>0.724</b>	0.673
	MU	0.722	<b>0.693</b>
<i>Ka-En</i>	MB	0.783	0.736
	MU	<b>0.824</b>	<b>0.778</b>
<i>Ma-En</i>	MB	<b>0.902</b>	<b>0.815</b>
	MU	0.899	<b>0.815</b>
<i>Mr</i>	MB	0.885	0.872
	MU	<b>0.895</b>	<b>0.887</b>
<i>Ta-En</i>	MB	0.829	0.791
	MU	<b>0.839</b>	<b>0.798</b>
<i>Ur-En</i>	MB	0.800	0.794
	MU	<b>0.816</b>	<b>0.810</b>
<i>Ur</i>	MB	<b>0.902</b>	<b>0.902</b>
	MU	0.895	0.895
<i>En</i>	MB	0.917	<b>0.917</b>
	MU	<b>0.918</b>	<b>0.917</b>

Table 2: Performance of ELFI using MB and MU. (Best performance is highlighted using bold).

not have any training instances. We name the language model as *AllBOne*.

**4.2.4 Instance transfer.** Here we go the other way, i.e., instead of directly evaluating the model in the target language, we first translate the target language instances to the source language using Google Translate API<sup>13</sup> on which the model is initially fine-tuned. In the zero shot setting **Ix0**, no gold target instances are used. In this scenario, we can again use a few gold target language instances (by translating to source language) to fine-tune the source model further. We name them **Ix32** and **Ix64**. As for the code-mixed instances, the translation across languages will be inaccurate; we limit this experiment to the monolingual setting only.

**4.2.5 Synthetic transfer.** Due to the less availability of low resource languages, in this technique, we experiment, can resource-rich language be useful if we translate them to low-resource language and build the model from scratch. To accomplish this, we translate (plentiful) gold source language (e.g., English) into the silver target language (e.g., Bengali, Hindi, etc.) instances using Google Translate API to train the model in the target language. This form of translation is widely used for model transfer; see Kozhevnikov and Titov [20]. Here, again, we can throw in zero or a few target gold instances, guiding to variations we name **STx0**, **STx32** and **STx64**.

### 4.3 Experimental setup

All models are assessed using the same 70:10:20 train, validation, and test split, stratified by class across the splits. For the transfer learning experiments, we use 32 and 64 training instances from each

<sup>13</sup><https://cloud.google.com/translate>

Lang	Model	Accuracy	M-F1
<i>Bn</i>	MB	<b>0.906</b> (+0.003)	<b>0.896</b> (+0.004)
	MU	<b>0.906</b> (+0.001)	0.895 (+0.001)
<i>Hi</i>	MB	0.799 (+0.023)	0.783 (+0.015)
	MU	<b>0.804</b> (-0.002)	<b>0.794</b> (-0.005)
<i>Hi-En</i>	MB	<b>0.692</b> (-0.032)	0.605 (-0.068)
	MU	0.636 (-0.086)	<b>0.622</b> (-0.071)
<i>Ka-En</i>	MB	0.803 (+0.02)	0.720 (-0.016)
	MU	<b>0.836</b> (+0.012)	<b>0.771</b> (-0.007)
<i>Ma-En</i>	MB	0.891 (-0.011)	0.797 (-0.018)
	MU	<b>0.894</b> (-0.005)	<b>0.802</b> (-0.013)
<i>Mr</i>	MB	0.889 (+0.004)	0.877 (+0.005)
	MU	<b>0.909</b> (+0.014)	<b>0.901</b> (+0.014)
<i>Ta-En</i>	MB	0.836 (+0.007)	0.786 (-0.005)
	MU	<b>0.845</b> (+0.006)	<b>0.795</b> (-0.003)
<i>Ur-En</i>	MB	0.763 (-0.037)	0.754 (-0.04)
	MU	<b>0.772</b> (-0.044)	<b>0.763</b> (-0.047)
<i>Ur</i>	MB	0.905 (+0.003)	0.905 (+0.003)
	MU	<b>0.906</b> (+0.011)	<b>0.906</b> (+0.011)
<i>En</i>	MB	<b>0.915</b> (-0.002)	<b>0.915</b> (-0.002)
	MU	0.912 (-0.006)	0.912 (-0.005)

Table 3: Performance of joint training. Best performance is highlighted using bold and values within parenthesis denote improvements over to ELFI.

class to further fine-tune the model in other languages. We make three such different random sets for each target dataset to make our evaluation more effective and report the average performance.

All models were coded in Python, using the Pytorch library. The models were run for 10 epoch with Adam optimizer, batch\_size = 16, learning\_rate =  $2e-5$  and adam\_epsilon =  $1e-8$ . We set the number of tokens  $n = 128$  for the experiments. We did not perform any hyper-parameter searching due to the limitation of computational resources; besides, the stated hyper-parameters have shown state-of-the-art performance in some of the previous shared-task [5, 10]. All our results are reported on the test set. For all the experiments, we use Ryzen 9, 5<sup>th</sup> gen 12 core processor, a Linux-based system with 64GB RAM and 16GB RTX 3080 GPU.

### 4.4 Evaluation metric

To remain consistent with existing literature, we assess our models in terms of **accuracy** and **macro F1-score**. These metrics together should be able to thoroughly evaluate the classification performance in distinguishing between abusive and normal posts. For zero-shot and few-shot settings, we report only **macro F1-score** due to scarcity of space.

## 5 RESULTS

### 5.1 Performance of ELFI

In Table 2, we show the performance of each model for all the languages in terms of accuracy and macro F1 score. We observe model’s performance varies from language to language. We see out of 10 language types, MuRIL (MU) outperforms m-BERT (MB) in

Lang	Model	Bn	Hi	Hi-En	Ka-En	Ma-En	Mr	Ta-En	Ur-En	Ur	En	AllBOne
Bn	MB	-	0.670	0.431	<b>0.641</b>	0.482	0.631	0.617	0.368	0.461	0.403	<b>0.678</b>
	MU	-	0.713	<b>0.497</b>	0.635	<b>0.633</b>	<b>0.698</b>	<b>0.625</b>	<b>0.415</b>	<b>0.763</b>	<b>0.510</b>	0.652
Hi	MB	0.391	-	<b>0.443</b>	0.573	0.434	0.514	0.499	<b>0.335</b>	0.404	0.386	0.555
	MU	<b>0.585</b>	-	0.395	<b>0.644</b>	<b>0.629</b>	<b>0.732</b>	<b>0.673</b>	0.323	<b>0.701</b>	<b>0.533</b>	<b>0.592</b>
Hi-En	MB	<b>0.444</b>	<b>0.548</b>	-	0.323	0.470	<b>0.501</b>	0.420	0.360	0.395	<b>0.464</b>	<b>0.382</b>
	MU	0.425	0.494	-	<b>0.326</b>	<b>0.528</b>	0.481	<b>0.508</b>	<b>0.389</b>	<b>0.477</b>	0.438	0.379
Ka-En	MB	<b>0.500</b>	<b>0.613</b>	0.503	-	0.536	0.505	<b>0.619</b>	0.366	0.467	0.439	0.586
	MU	0.455	0.609	<b>0.550</b>	-	<b>0.649</b>	<b>0.583</b>	<b>0.702</b>	<b>0.391</b>	<b>0.537</b>	<b>0.442</b>	<b>0.628</b>
Ma-En	MB	0.572	<b>0.545</b>	0.534	0.485	-	<b>0.552</b>	<b>0.627</b>	0.292	0.493	<b>0.466</b>	0.569
	MU	0.499	0.496	<b>0.559</b>	<b>0.500</b>	-	0.511	0.604	<b>0.340</b>	<b>0.563</b>	0.464	<b>0.603</b>
Mr	MB	0.406	<b>0.661</b>	0.400	0.580	0.430	-	<b>0.606</b>	<b>0.523</b>	0.472	0.394	0.599
	MU	<b>0.552</b>	<b>0.813</b>	<b>0.434</b>	<b>0.681</b>	0.595	-	<b>0.703</b>	0.403	<b>0.710</b>	<b>0.468</b>	<b>0.816</b>
Ta-En	MB	<b>0.511</b>	<b>0.592</b>	0.509	0.593	<b>0.614</b>	0.564	-	0.340	0.450	0.430	<b>0.614</b>
	MU	0.459	0.516	<b>0.567</b>	<b>0.634</b>	<b>0.621</b>	<b>0.572</b>	-	<b>0.355</b>	<b>0.513</b>	<b>0.442</b>	<b>0.627</b>
Ur-En	MB	<b>0.326</b>	<b>0.382</b>	0.315	0.467	0.404	<b>0.468</b>	<b>0.484</b>	-	0.298	0.296	<b>0.359</b>
	MU	0.317	0.322	<b>0.354</b>	<b>0.473</b>	<b>0.420</b>	0.399	0.464	-	<b>0.347</b>	<b>0.302</b>	0.340
Ur	MB	0.377	<b>0.655</b>	<b>0.357</b>	0.597	0.426	0.548	0.652	0.404	-	0.342	<b>0.678</b>
	MU	<b>0.580</b>	<b>0.804</b>	0.339	<b>0.652</b>	<b>0.635</b>	<b>0.803</b>	<b>0.669</b>	<b>0.453</b>	-	<b>0.462</b>	<b>0.753</b>
En	MB	0.352	0.626	<b>0.481</b>	0.577	0.398	0.397	0.528	0.324	0.350	-	<b>0.668</b>
	MU	<b>0.606</b>	<b>0.826</b>	0.476	<b>0.656</b>	<b>0.443</b>	<b>0.759</b>	<b>0.638</b>	<b>0.354</b>	<b>0.729</b>	-	<b>0.732</b>

Table 4: Macro F1 score of MTx0. Row-wise language names represent the target language, and column-wise language names represent the source language model. The better among MB and MU is highlighted in bold. Blue denotes the best source + model pair for a target language. Blue represents the best, and green the second-best performing model per row (MU or MB).

Lang	Model	Shot	Bn	Hi	Hi-En	Ka-En	Ma-En	Mr	Ta-En	Ur-En	Ur	En	AllBOne
Bn	MB	MTx32	-	0.705	<b>0.647</b>	0.686	0.664	0.726	0.705	0.579	<b>0.739</b>	0.71	0.725
		MTx64	-	(+0.055)	(+0.033)	(+0.034)	(+0.046)	(+0.024)	(+0.045)	(+0.071)	(+0.011)	(+0.04)	(+0.015)
	MU	MTx32	-	<b>0.777</b>	0.579	<b>0.725</b>	<b>0.712</b>	<b>0.76</b>	<b>0.718</b>	<b>0.628</b>	<b>0.777</b>	<b>0.755</b>	<b>0.774</b>
		MTx64	-	(+0.003)	(+0.131)	(+0.015)	(+0.028)	(+0.01)	(+0.032)	(+0.082)	(+0.003)	(+0.015)	(-0.004)
Hi	MB	MTx32	0.656	-	<b>0.613</b>	<b>0.679</b>	0.465	0.672	0.647	0.569	<b>0.677</b>	0.643	0.635
		MTx64	(+0.014)	-	(+0.027)	(+0.011)	(+0.115)	(+0.008)	(+0.033)	(+0.031)	(+0.013)	(+0.037)	(+0.025)
	MU	MTx32	<b>0.739</b>	-	0.607	<b>0.696</b>	<b>0.678</b>	<b>0.732</b>	<b>0.702</b>	<b>0.624</b>	<b>0.737</b>	<b>0.751</b>	<b>0.755</b>
		MTx64	(+0.011)	-	(+0.063)	(+0.014)	(+0.032)	(+0.028)	(+0.018)	(+0.056)	(+0.013)	(+0.009)	(+0.015)
Hi-En	MB	MTx32	<b>0.538</b>	0.552	-	0.514	0.397	<b>0.531</b>	0.538	0.487	<b>0.544</b>	<b>0.567</b>	0.485
		MTx64	(+0.012)	(+0.018)	-	(+0.026)	(+0.133)	(+0.019)	(+0.032)	(+0.023)	(+0.006)	(+0.013)	(+0.045)
	MU	MTx32	0.521	<b>0.553</b>	-	<b>0.575</b>	<b>0.431</b>	0.53	<b>0.552</b>	<b>0.496</b>	0.532	<b>0.581</b>	<b>0.508</b>
		MTx64	(+0.029)	(+0.027)	-	(+0.055)	(-0.051)	(+0.02)	(+0.038)	(+0.054)	(+0.028)	(+0.019)	(+0.042)
Ka-En	MB	MTx32	0.589	0.604	<b>0.585</b>	-	0.606	0.57	0.598	0.53	0.584	0.61	<b>0.624</b>
		MTx64	(+0.031)	(+0.036)	(+0.045)	-	(+0.034)	(+0.06)	(+0.082)	(+0.06)	(+0.026)	(+0.03)	(+0.036)
	MU	MTx32	<b>0.618</b>	<b>0.638</b>	0.551	-	<b>0.615</b>	<b>0.572</b>	<b>0.664</b>	<b>0.561</b>	<b>0.614</b>	<b>0.62</b>	<b>0.676</b>
		MTx64	(+0.032)	(+0.012)	(+0.069)	-	(+0.005)	(+0.058)	(+0.036)	(+0.029)	(+0.036)	(+0.02)	(+0.004)
Ma-En	MB	MTx32	0.597	0.562	<b>0.566</b>	0.58	-	<b>0.58</b>	<b>0.643</b>	0.484	0.567	0.582	<b>0.636</b>
		MTx64	(+0.033)	(+0.068)	(+0.054)	(+0.05)	-	(+0.04)	(+0.017)	(+0.066)	(+0.063)	(+0.038)	(+0.014)
	MU	MTx32	<b>0.6</b>	<b>0.6</b>	0.537	<b>0.618</b>	-	0.555	<b>0.624</b>	<b>0.525</b>	<b>0.627</b>	<b>0.589</b>	0.604
		MTx64	(+0.04)	(+0.04)	(+0.083)	(+0.032)	-	(+0.035)	(+0.036)	(+0.045)	(+0.023)	(+0.031)	(+0.036)
Mr	MB	MTx32	0.73	<b>0.762</b>	<b>0.629</b>	0.725	0.515	-	0.74	0.617	0.729	0.708	<b>0.757</b>
		MTx64	(+0.03)	(+0.038)	(+0.041)	(+0.035)	(+0.225)	-	(+0.05)	(+0.023)	(+0.041)	(+0.062)	(+0.023)
	MU	MTx32	<b>0.784</b>	<b>0.834</b>	0.608	<b>0.753</b>	<b>0.725</b>	-	<b>0.795</b>	<b>0.701</b>	<b>0.812</b>	<b>0.801</b>	<b>0.839</b>
		MTx64	(+0.026)	(+0.016)	(+0.102)	(+0.027)	(+0.065)	-	(+0.035)	(+0.039)	(-0.002)	(+0.019)	(+0.001)
Ta-En	MB	MTx32	0.608	0.619	0.587	<b>0.631</b>	<b>0.622</b>	0.607	-	0.537	0.606	0.616	<b>0.661</b>
		MTx64	(+0.012)	(+0.001)	(+0.013)	(+0.009)	(-0.002)	(+0.003)	-	(+0.033)	(+0.004)	(+0.004)	(-0.001)
	MU	MTx32	<b>0.628</b>	<b>0.648</b>	<b>0.592</b>	<b>0.654</b>	0.61	<b>0.621</b>	-	<b>0.547</b>	<b>0.622</b>	<b>0.619</b>	<b>0.673</b>
		MTx64	(+0.002)	(+0.012)	(+0.028)	(+0.016)	(+0)	(+0.019)	-	(+0.043)	(+0.018)	(+0.011)	(+0.007)
Ur-En	MB	MTx32	<b>0.518</b>	0.489	0.486	<b>0.535</b>	0.461	<b>0.559</b>	<b>0.514</b>	-	<b>0.55</b>	<b>0.516</b>	<b>0.475</b>
		MTx64	(+0.022)	(+0.021)	(+0.034)	(+0.025)	(+0.009)	(+0.011)	(+0.016)	-	(+0.01)	(+0.014)	(+0.005)
	MU	MTx32	0.454	<b>0.5</b>	<b>0.491</b>	0.501	<b>0.506</b>	0.461	<b>0.502</b>	-	0.458	0.492	0.448
		MTx64	(+0.016)	(+0.03)	(+0.009)	(+0.009)	(+0.024)	(+0.059)	(+0.018)	-	(+0.042)	(+0.018)	(+0.012)
Ur	MB	MTx32	0.785	0.782	0.669	<b>0.75</b>	0.641	<b>0.802</b>	<b>0.803</b>	0.623	-	0.769	0.793
		MTx64	(+0.055)	(+0.038)	(+0.061)	(+0.05)	(+0.109)	(+0.028)	(+0.027)	(+0.057)	-	(+0.051)	(+0.047)
	MU	MTx32	<b>0.839</b>	<b>0.83</b>	<b>0.686</b>	0.71	<b>0.776</b>	<b>0.808</b>	0.756	<b>0.663</b>	-	<b>0.811</b>	<b>0.824</b>
		MTx64	(+0.011)	(+0.01)	(+0.024)	(+0.05)	(+0.034)	(+0.012)	(+0.054)	(+0.067)	-	(+0.019)	(+0.006)
En	MB	MTx32	0.768	<b>0.825</b>	<b>0.763</b>	<b>0.78</b>	0.695	<b>0.811</b>	<b>0.791</b>	<b>0.686</b>	0.793	-	0.798
		MTx64	(+0.042)	(+0.015)	(+0.057)	(+0.04)	(+0.095)	(+0.019)	(+0.029)	(+0.054)	(+0.037)	-	(+0.022)
	MU	MTx32	<b>0.815</b>	<b>0.837</b>	0.656	0.754	<b>0.727</b>	<b>0.829</b>	0.79	0.633	<b>0.824</b>	-	<b>0.801</b>
		MTx64	(+0.015)	(+0.013)	(+0.074)	(+0.026)	(+0.073)	(+0.011)	(+0.02)	(+0.117)	(+0.016)	-	(+0.029)

Table 5: Macro F1 score of MTx{32, 64} for different source and target languages. Row-wise language names represent the target language, and column-wise language names represent the source language model. The better among MB and MU is highlighted in bold. Blue denotes the best source + model pair for a target language. Blue represents the best, and green the second-best performing model per row (MU or MB). For MTx64, we show the gain added over MTx32.

7 languages in terms of macro F1-score. Although **MU** was pre-trained in only 17 languages (quite less compared to **MB**), it was specifically focused on Indian languages. We posit that this is one of the primary reasons for its prevailing superior performance over **MB**.

## 5.2 Performance of joint training/cross-lingual training

Here we investigate the importance of joint training. Even though different language has different character symbols, the main intention for creating all these datasets have been to identify abusive post. Therefore, in order to check if the dataset in one language can be useful for another language or not, we execute joint training by merging all the training language instances. In Table 3 we summarize our results. We observe for some languages, the performance improved, and for some, the performance decreased. In general, **MU** works better than **MB** as observed in case of ELFI. Furthermore, we notice the performance deviation is negligible compared to the ELFI setting. Thus, as joint training is usually more expensive in terms of computational resource usage, it is more reasonable to perform self-training in a resource-constrained environment.

## 5.3 Performance of model transfer

**5.3.1 Performance of MTx0.** In Table 4, we show the performance of zero-shot model transfer outcomes across all pairs of languages. As expected, the macro F1 scores are worse compared to the ELFI setting. Further **MU** works better than **MB**.

- **Ur + MU** is the most effective model when the target language is **Bn**; **Hi + MU** is the next most effective one. If **Hi** is the target language, **Mr + MU** is the most effective model followed by **Ur + MU**. Likewise, when **Ur** is the target language, **Hi + MU** model performs the best followed **Mr + MU**. If **Mr** is the target language although the **AllBOne + MU** model performs the best, **Hi + MU** is the closest second. This may be explained by their membership in the Indo-Aryan language family.
- The **Hi + MB** model performs the best for the **Hi-En** target language. This is possibly because of the fact that **Hi-En** has the same underlying semantics as that of **Hi**.
- **Ta-En + MU** model is most effective for the **Ka-En** language. Conversely **Ka-En + MU** is the best source for the **Ta-En** language. Besides, for the target language **Ma-En**, **Ta-En + MB** model obtains the highest performance. The effectiveness of these models is possibly because of the fact that all these languages belong to the Dravidian language family.

Overall we observe that zero-shot abusive language detection can be useful when the source and the target languages are from the same language family.

**5.3.2 Performance of MTx32 and MTx64.** Table 5 illustrates the effect of allowing a small portion of available gold instances in the target language for second-stage fine-tuning. We observe adding more instances improves the overall performance in both **MU** and **MB** settings. With more target language gold instances, the deficit from ELFI continues to decrease. An additional interesting

Lang	Model	Bn	Hi	Mr	Ur	En
<b>Bn</b>	<b>MB</b>	-	<b>0.677</b>	<b>0.59</b>	0.635	<b>0.605</b>
	<b>MU</b>	-	0.621	0.539	<b>0.666</b>	0.595
<b>Hi</b>	<b>MB</b>	<b>0.588</b>	-	0.519	0.588	<b>0.597</b>
	<b>MU</b>	0.53	-	<b>0.658</b>	<b>0.614</b>	0.571
<b>Mr</b>	<b>MB</b>	<b>0.487</b>	<b>0.617</b>	-	0.53	0.511
	<b>MU</b>	0.445	<b>0.609</b>	-	<b>0.568</b>	<b>0.519</b>
<b>Ur</b>	<b>MB</b>	<b>0.563</b>	<b>0.627</b>	0.529	-	<b>0.524</b>
	<b>MU</b>	0.498	<b>0.601</b>	<b>0.552</b>	-	0.507

**Table 6: Macro F1 score of Ix0. Row-wise language names represent the target language, and column-wise language names represent the source language model. The better among MB and MU is highlighted in bold. Blue denotes the best source + model pair for a target language. Blue represents the best, and green the second-best performing model per row (MU or MB).**

observation is that the combinations that were best in the zero-shot setting get altered in many cases here.

- Here, both **Ur + MU** and **Hi + MU** models perform best for target language **Bn**. Conversely, **Bn + MU** is the best source for the target language **Ur**.
- For the target languages **Hi**, **Mr**, **Ka-En**, and **Ta-En**, we observe **AllBOne + MU** model is the most effective. This could be due to the advantages drawn by this model from the two stage fine-tuning. While in the first stage it learns the diversity associated with the languages from different linguistic families, in the second stage it learns the target language specific intricacies from the few-shot labels provided.
- Further we notice that for the target language **Hi-En**, **En + MU** model performs the best. This may be due to the **Hi-En** examples shown to the **En + MU** model as few-shots from which it is able to learn the English-Hindi switching patterns and their semantic connections.
- Likewise zero-shot performance, we observe that the **Ta-En + MB** model achieves the best score on **Ma-En**.
- The **Mr + MB** model is most effective for the target language **Ur-En**, followed by the **Ur** model.
- For the target language **En**, similar to zero-shot performance, **Hi + MU** is the best source model.

Overall we observe that in a few shot setting, the **AllBOne** model is the most effective.

## 5.4 Performance of instance transfer

In this section, we compare the performance of instance transfer (**Ix**) with model transfer (**MTx**). Table 6 shows the performance of **Ix0**. We observe, unlike **MTx**, **MB** performs relatively better than **MU**. Although comparing Table 4 and 6, we notice **MTx0** outperforms **Ix0** for all the languages. We observe two reasons for the inferior performance of **Ix0** – (a) while translating one language to another, translation errors are inevitable and (b) translation tones down the level abusiveness of a post.

Lang	Model	Shot	<i>Bn</i>	<i>Hi</i>	<i>Mr</i>	<i>Ur</i>	<i>En</i>
<i>Bn</i>	MB	Ix32	-	0.679	0.64	0.695	0.679
		Ix64	-	(+0.011)	(+0.02)	(+0.015)	(+0.031)
	MU	Ix32	-	0.716	0.695	0.724	0.704
		Ix64	-	(+0.004)	(+0.015)	(+0.006)	(+0.016)
<i>Hi</i>	MB	Ix32	0.661	-	0.691	0.66	0.685
		Ix64	(+0.009)	-	(+0.009)	(+0.01)	(+0.015)
	MU	Ix32	0.697	-	0.732	0.717	0.719
		Ix64	(+0.023)	-	(+0.008)	(+0.013)	(+0.001)
<i>Mr</i>	MB	Ix32	0.624	0.678	-	0.653	0.689
		Ix64	(+0.046)	(+0.002)	-	(+0.027)	(+0.021)
	MU	Ix32	0.683	0.742	-	0.732	0.696
		Ix64	(+0.017)	(-0.002)	-	(+0.008)	(+0.014)
<i>Ur</i>	MB	Ix32	0.692	0.685	0.659	-	0.652
		Ix64	(+0.038)	(+0.025)	(+0.021)	-	(+0.038)
	MU	Ix32	0.722	0.706	0.674	-	0.644
		Ix64	(+0.018)	(+0.024)	(+0.016)	-	(+0.046)

Table 7: Macro F1 score of Ix{32, 64} for different source and target languages. Row-wise language names represent the target language, and column-wise language names represent the source language model. The better among MB and MU is highlighted in bold. Blue denotes the best source + model pair for a target language. Blue represents the best, and green the second-best performing model per row (MU or MB). For Ix64, we show the gain over Ix32.

Lang	Model	STx0	STx32	STx64
<i>Bn</i>	MB	0.572	0.726	0.761
	MU	0.612	0.739	0.779
<i>Hi</i>	MB	0.601	0.672	0.705
	MU	0.617	0.666	0.674
<i>Mr</i>	MB	0.725	0.779	0.800
	MU	0.698	0.802	0.807
<i>Ur</i>	MB	0.654	0.774	0.821
	MU	0.664	0.810	0.833

Table 8: Macro F1 for STx{0,32,64}. Best performance is highlighted in bold.

Further, we also experiment with the few-shot setup in instance transfer environment (i.e., Ix16 and Ix32). Table 7 shows the performance, which although is better than Ix0, cannot outperform the numbers obtained from the model transfer schemes (see Table 5 vs. Table 7). Overall we maintain that for abusive language detection, model transfer schemes are superior to instance transfer schemes.

### 5.5 Performance of synthetic transfer

Here we explore the significance of silver instances in a low-resource setting. In the scarcity of ample target instances for ELFI style training, we can leverage potentially plentiful instances from a source language, i.e., *En*. Table 8 shows the results. Even after using all available silver target instances, we observe that adding gold target instances gives steady improvements. This implies the worth of always having at least some gold target instances in such tasks.

## 6 ERROR ANALYSIS

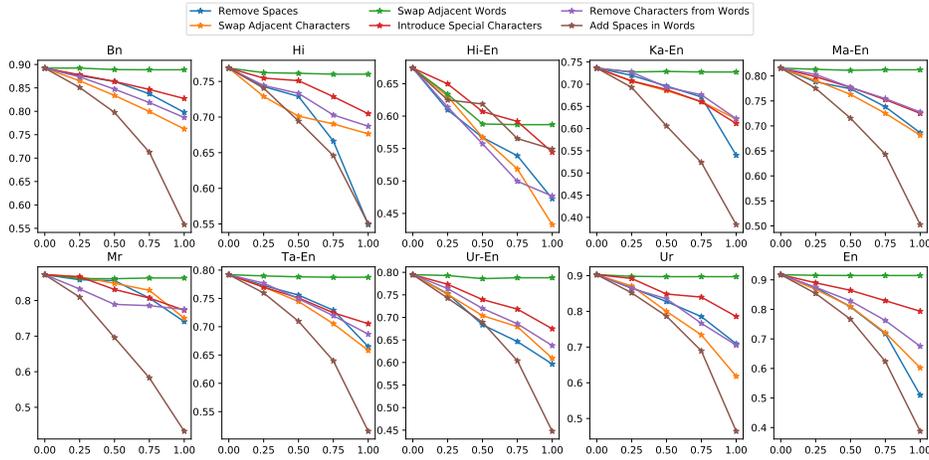
To delve deeper, we conduct an error analysis on both models using a small number of test data points wrongly classified by the models.

We examine our error from three independent directions – **general error**, **instance transfer error**, and **synthetic transfer error**.

**General error:** We analyze the common errors and segregate them into the following four categories.

- The **presence of discriminatory features** is one of the crucial points of consideration for labeling a post abusive. We observe that models misclassify such posts where some specific words are present in the post that can be used in both abusive and non-abusive contexts. For example, the word ‘nigga’ in a sentence is more likely to be considered abusive, whereas African-American people use the word in a non-abusive context. For instance, in the following tweet, “Niggas pulled up on a nigga said hey come to HR for a drug test. Homie said word bet I quit LMFAOOOO” is predicted as abusive by the model, but its actual label is normal. On the other hand “I play wit pussy not these niggas” is correctly predicted abusive by the model. We see m-BERT suffers mostly for this type of error.
- If the **contextual information** is limited in a post, models cannot predict its actual label. We observe some posts with no abusive text; however, its context makes it abusive possibly because of the presence of emoji, or attached URL. For example, the Bengali instance “Ai maa\*i tui holo akta 🍆🍆🍆🍆🍆<sup>14</sup>” (Translation: You woman, you are a 🍆🍆🍆🍆🍆) without the emoji can be considered as non abusive; but the presence of the emojis makes it abusive. Our analysis identified that MU suffers mostly for this category of errors.
- **Implicit abusiveness** is another form of an error where the model fails to capture the context due to the absence of explicit abusive content. Understanding these instances requires understanding of sarcasm, complicated reasoning

<sup>14</sup>We write non-English posts using Roman scripts.



**Figure 1: Macro F1 score under the adversarial attack. The X-axis represents (%) of attack; Y-axis represents changes in macro F1 score.**

skills, or background knowledge of some situations. For instance the Hindi post “Khujliwal ko karbwakar ashabaadi bana do <sup>14</sup>” (Translation: Make Khujliwal optimistic by getting it done.) is abusive toward a politician (Kejriwal) in India; nevertheless, the model fails to predict its actual label as it is sarcastic and need background knowledge about the politician. We notice **MB** suffers mostly for this category of errors.

- **Tentatively disputed annotations** are another type of situation where the ground truth label is ambiguous. The context of a post can be multi-dimensional, and based on annotators’ understanding, the labeling can be biased toward a specific direction. For example, the following Bengali post “Kon kon bokach\*da mile man of the match nirbachan korche. <sup>14</sup>” (Translation: Which idiot-fucker is choosing Man of the Match.) has been wrongly annotated as normal, though the use of slur words makes it a clearly abusive post.

**Instance transfer error:** In the case of instance transfer, where we translate the target language model to the source language model, we observe the translation affects the true semantics of the actual target posts. Usually, the translator tool reduces the toxicity score of a sentence while translating; sometimes, it cannot translate appropriately because of context-sensitive situations. Hence model performance deteriorates while predicting those instances. For example, “Tui akta kukurer bachha <sup>14</sup>” (a slur in Bengali) has been translated to “You’re a puppy” and does not look toxic.

**Synthetic transfer error:** A similar observation holds for synthetic transfer experiments. Here, the model is trained primarily on relatively less intense words. Hence, different kinds of abusive words/phrases get translated to either the same word or acquire some other representations. Naturally therefore the model performance degrades. For instance, when the following sentence “Ladies and Gentlemen, Victoria Soliz has fucked up once again” is translated to Hindi, it becomes “Deviyon aur sajjanon, Viktoriya Solis nek baar phir gadbad kar di hai <sup>14</sup>” (Translation: Ladies and gentlemen, Victoria Solies has messed up once again). Though the original

post was abusive, the translation made the sentence less intense and could be considered as normal.

## 7 ROBUSTNESS AGAINST ADVERSARIAL CHANGES IN THE TEXT

Observing the exceptional performance of ELFI style training, in this section, we try to understand, even if we can build an immaculate abusive speech classification model, how robust these models are against various adversarial changes. Specifically, we propose the following six black-box ‘attacks’, which assumes that the attackers do not know the details of the detection algorithm, hence only creating the best possible guess to avoid exposure.

*Remove spaces:* The spaces between adjacent words are removed. As the word-based language models primarily rely on tokenizing the text into a sequence of words, removing all spaces leads to a single <unk> token. Thus, finding the token boundaries would be difficult for the model.

*Add spaces in words:* Conversely, we introduce spaces within individual words in a text, making the word unrecognizable based on characters treated as word separators.

*Remove characters from words:* Here we remove the characters from words to introduce typos. These words will still be readable by a human, but the tokens are changed from the classifier’s perspective.

*Introduce special characters:* Here we introduce random special characters within words to introduce typos.

*Swap adjacent characters:* We introduce another form of typo within words by swapping adjacent characters.

*Swap adjacent words:* The transformer-based model tries to understand the underlying meaning of a sentence based on the way words are ordered. So, we swap adjacent words to change the relative ordering of words.

**Observations:** In Figure 1, we demonstrate the performance of all the language models under various adversarial attacks for **MB**<sup>15</sup>. We see the performance of all the language models drops gradually

<sup>15</sup>Similar observations hold for **MU** as well.

with the increasing amount of adversarial attacks. Except for *Hi-En*, all other models exhibit robustness against **adjacent word swap**, which indicates that the presence of abusive words is itself a strong signal for the model to judge whether a text is abusive or not irrespective of the relative order of the words. All the language models suffer when spaces are introduced within words, making it difficult to understand the actual word boundaries for the model. Overall similar observations are found for other types of adversarial attacks, which opens up another dimension of work that needs to be addressed to strengthen the models against such adversarial attacks.

## 8 CONCLUSION

In this paper, we tried to perform multilingual abusive language detection for Indic languages. We used transformer-based models to develop classifiers for abusive speech identification, using eight different languages from 14 publicly available resources. We perform several experiments for multiple languages under various settings - ELFI, zero-shot learning, few-shot learning, model transfer, instance transfer, cross-lingual learning, etc. Overall, we noticed that model transfer obtains better performance than instance transfer; further model transfer is advantageous when the source model language and target language belong to the same language family. Further in a few-shot setting the *AILBoNe* model performs the best as it gains from both the fine-tuning stages; while in the first stage it learns universal features, in the second stage it learns language specific features. We observed for low-resource languages, though synthetic silver instances are helpful to build classifiers for abusive language detection, further fine-tuning the model with gold target instances shows steady improvements.

One of the main drawbacks we observed was that the model's performance decreased against adversarial attacks. We plan to improve the existing models to make them agnostic of the adversarial attacks. Further, we plan to create datasets in other regional languages using the knowledge we obtained from our experiments.

## REFERENCES

- [1] 2017. Moderators who had to view Child abuse content sue microsoft, claiming PTSD. <https://www.theguardian.com/technology/2017/jan/11/microsoft-employees-child-abuse-lawsuit-ptsd>
- [2] Muhammad Pervez Akhter, Zheng Jiangbin, Irfan Raza Naqvi, Mohammed Abdelmajeed, and Muhammad Tariq Sadiq. 2020. Automatic detection of offensive language for urdu and roman urdu. *IEEE Access* 8 (2020), 91213–91226.
- [3] Maaz Amjad, Alisa Zhila, Grigori Sidorov, Andrey Labunets, Sabur Butt, Hamza Imam Amjad, Oxana Vitman, and Alexander Gelbukh. 2021. UrduThreat@ FIRE2021: Shared Track on abusive threat Identification in Urdu. In *FIRE 2021: Forum for Information Retrieval Evaluation, December 13-17, 2021, India*. ACM.
- [4] Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, and Vasudeva Varma. 2017. Deep learning for hate speech detection in tweets. In *Proceedings of the 26th international conference on World Wide Web companion*. 759–760.
- [5] Somnath Banerjee, Maulindu Sarkar, Nancy Agrawal, Punyajoy Saha, and Mithun Das. 2021. Exploring Transformer Based Models to Identify Hate Speech and Offensive Content in English and Indo-Aryan Languages. *arXiv preprint arXiv:2111.13974* (2021).
- [6] Aditya Bohra, Deepanshu Vijay, Vinay Singh, Syed Sarfaraz Akhtar, and Manish Shrivastava. 2018. 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*. 36–41.
- [7] Bharathi Raja Chakravarthi, Ruba Priyadarshini, Navya Jose, Thomas Mandl, Prasanna Kumar Kumaresan, Rahul Ponnusamy, RL Hariharan, John Philip McCrae, Elizabeth Sherly, et al. 2021. Findings of the shared task on offensive language identification in Tamil, Malayalam, and Kannada. In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*. 133–145.
- [8] Ying Chen, Yilu Zhou, Sencun Zhu, and Heng Xu. 2012. Detecting offensive language in social media to protect adolescent online safety. In *2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Conference on Social Computing*. IEEE, 71–80.
- [9] Gokul Chittaranjan, Yogarshi Vyas, Kalika Bali, and Monojit Choudhury. 2014. Word-level language identification using crf: Code-switching shared task report of msr india system. In *Proceedings of The First Workshop on Computational Approaches to Code Switching*. 73–79.
- [10] Mithun Das, Somnath Banerjee, and Punyajoy Saha. 2021. Abusive and Threatening Language Detection in Urdu using Boosting based and BERT based models: A Comparative Approach. *arXiv preprint arXiv:2111.14830* (2021).
- [11] Mithun Das, Binny Mathew, Punyajoy Saha, Pawan Goyal, and Animesh Mukherjee. 2020. Hate speech in online social media. *ACM SIGWEB Newsletter* Autumn (2020), 1–8.
- [12] Mithun Das, Punyajoy Saha, Ritam Dutt, Pawan Goyal, Animesh Mukherjee, and Binny Mathew. 2021. You too brutus! trapping hateful users in social media: Challenges, solutions & insights. In *Proceedings of the 32nd ACM Conference on Hypertext and Social Media*. 79–89.
- [13] Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 11.
- [14] J. Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*.
- [15] Antigoni Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. In *Twelfth International AAAI Conference on Web and Social Media*.
- [16] Saurabh Sampatrao Gaikwad, Tharindu Ranasinghe, Marcos Zampieri, and Christopher Homan. 2021. Cross-lingual Offensive Language Identification for Low Resource Languages: The Case of Marathi. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*. 437–443.
- [17] Nicola F Johnson, R Leahy, N Johnson Restrepo, Nicolas Velasquez, Ming Zheng, P Manrique, P Devkota, and Stefan Wuchty. 2019. Hidden resilience and adaptive dynamics of the global online hate ecology. *Nature* 573, 7773 (2019), 261–265.
- [18] Muhammad Moin Khan, Khurram Shahzad, and Muhammad Kamran Malik. 2021. Hate Speech Detection in Roman Urdu. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 20, 1 (2021), 1–19.
- [19] Simran Khanuja, Diksha Bansal, Sarvesh Mehtani, Savya Khosla, Atreyee Dey, Balaji Gopalan, Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja Nagipogu, Shachi Dave, et al. 2021. MuriL: Multilingual representations for indian languages. *arXiv preprint arXiv:2103.10730* (2021).
- [20] Mikhail Kozhevnikov and Ivan Titov. 2014. Cross-lingual model transfer using feature representation projection. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 579–585.
- [21] Ritesh Kumar, Atul Kr Ojha, Shervin Malmasi, and Marcos Zampieri. 2018. Benchmarking aggression identification in social media. In *Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018)*. 1–11.
- [22] Ritesh Kumar, Atul Kr Ojha, Marcos Zampieri, and Shervin Malmasi. 2018. Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018). In *Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018)*.
- [23] Thomas Mandl, Sandip Modha, Anand Kumar M, and Bharathi Raja Chakravarthi. 2020. Overview of the hasoc track at fire 2020: Hate speech and offensive language identification in tamil, malayalam, hindi, english and german. In *Forum for Information Retrieval Evaluation*. 29–32.
- [24] Thomas Mandl, Sandip Modha, Prasenjit Majumder, Daksh Patel, Mohana Dave, Chintak Mandlia, and Aditya Patel. 2019. Overview of the hasoc track at fire 2019: Hate speech and offensive content identification in indo-european languages. In *Proceedings of the 11th forum for information retrieval evaluation*. 14–17.
- [25] Thomas Mandl, Sandip Modha, Gautam Kishore Shahi, Hiren Madhu, Shrey Satapara, Prasenjit Majumder, Johannes Schaefer, Tharindu Ranasinghe, Marcos Zampieri, Durgesh Nandini, et al. 2021. Overview of the HASOC subtrack at FIRE 2021: Hate speech and offensive content identification in English and Indo-Aryan languages. *arXiv preprint arXiv:2112.09301* (2021).
- [26] Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. HateXplain: A Benchmark Dataset for Explainable Hate Speech Detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 14867–14875.
- [27] Arijit Nag, Bidisha Samanta, Animesh Mukherjee, Niloy Ganguly, and Soumen Chakrabarti. 2021. A Data Bootstrapping Recipe for Low-Resource Multilingual Relation Classification. In *Proceedings of the 25th Conference on Computational Natural Language Learning*. 575–587.
- [28] Casey Newton. 2019. The terror queue. <https://www.theverge.com/2019/12/16/21021005/google-youtube-moderators-ptsd-accenture-violent-disturbing-content-interviews-video>

- [29] Alexandra Olteanu, Carlos Castillo, Jeremy Boy, and Kush Varshney. 2018. The effect of extremist violence on hateful speech online. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 12.
- [30] Georgios K. Pitsilis, H. Ramampiaro, and H. Langseth. 2018. Detecting Offensive Language in Tweets Using Deep Learning. *ArXiv abs/1801.04433* (2018).
- [31] Tharindu Ranasinghe and Marcos Zampieri. 2020. Multilingual offensive language identification with cross-lingual embeddings. *arXiv preprint arXiv:2010.05324* (2020).
- [32] Tharindu Ranasinghe and Marcos Zampieri. 2021. An evaluation of multilingual offensive language identification methods for the languages of india. *Information* 12, 8 (2021), 306.
- [33] Hammad Rizwan, Muhammad Haroon Shakeel, and Asim Karim. 2020. Hate-speech and offensive language detection in roman Urdu. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2512–2522.
- [34] Nauros Romim, Mosahed Ahmed, Hriteshwar Talukder, and Md Saiful Islam. 2021. Hate speech detection in the bengali language: A dataset and its baseline evaluation. In *Proceedings of International Joint Conference on Advances in Computational Intelligence*. Springer, 457–468.
- [35] Wiktor Soral, Michał Bilewicz, and Mikołaj Winiewski. 2018. Exposure to hate speech increases prejudice through desensitization. *Aggressive behavior* 44, 2 (2018), 136–146.
- [36] N Statt. 2017. YouTube is facing a full-scale advertising boycott over hate speech. *The Verge* (2017).
- [37] Janikke Solstad Vedeler, Terje Olsen, and John Eriksen. 2019. Hate speech harms: a social justice discussion of disabled Norwegians' experiences. *Disability & Society* 34, 3 (2019), 368–383.
- [38] Zeerak Waseem and Dirk Hovy. 2016. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop*, 88–93.
- [39] Ziqi Zhang, David Robinson, and Jonathan Tepper. 2018. Detecting hate speech on twitter using a convolution-gru based deep neural network. In *European semantic web conference*. Springer, 745–760.