

Hate Speech Detection in Hindi language using BERT and Convolution Neural Network

Shubham Shukla¹, Sushama Nagpal², Sangeeta Sabharwal³

Department of Computer Science and Engineering

Netaji Subhas University of Technology

Sector-3 Dwarka, Delhi, India

shubhcs7272@gmail.com¹, sushma.nagpal@nsut.ac.in², ssab@nsut.ac.in³

Abstract— Social media has become crucial in our lives; it inculcates our opinions by providing untreated information. Whether we might be not participating actively but indirectly everyone became part of its coverage. Wide spread of information over the internet without any validation made it hard to analyze the impact of misleading information. Cyber hate, which is used as a tool to incite violence against a group of people based on ethnicity, nationality, language, sexual orientation, religious faiths, etc., poses a disgraceful utilization of social media. Previous apposite studies reported hate speech mainly in the English language. Less effort has been made for the resource-constraint language such as Hindi, Marathi, Kannada, etc. This work entitles hate speech detection in low-resource Hindi language using BERT and Deep Convolution Neural Network. The proposed Hindi Hate Speech BERT Convolution Neural Network model intends to detect hate speech in real-time so that any harmful incidence can be avoided as early as possible. This model presents a two-stage architecture: In the first stage, we have applied a pre-trained BERT encoder to generate encodings. In the second stage, a convolution neural network followed by a sigmoid layer is used to detect text as hatred or non-hatred. Our model achieved 0.84 & 0.77 f1-score for Hasoc 2020 and Hasoc 2021 dataset respectively.

Keywords— Hate speech detection; BERT; convolution neural network; social media; Hindi text;

I. INTRODUCTION

While the Internet provides several potential avenues for accessing information, training, and communication, it also presents new risks, threats, and injuries. In particular, these risks and threats are most evident on Social Media Platforms (SMPs) like Facebook, Twitter, Tumblr, etc. Many of these platforms are free, and as such, users have the ability to access and share offensive and hateful content. However, an effective way to address the problem of misinformation on social media is through fact-checkers, who can provide context and background information about a particular topic or claim that can be critical in the context of a rapidly-spreading false narrative. However, fact-checkers are problematic because they are often biased and can be used to target political opponents. So, one of the most effective ways to address this problem is to use a machine learning based automated hate speech detection system [1]–[4]. Automated hate speech detection is a technique used to identify potentially offensive or hateful content in a digital context, such as a social media post. As such, hate speech detection is a more neutral approach that can identify

potentially offensive or hateful content without human interpretation or bias. Although hate speech detection and fact-checkers work efficiently in addressing the problem of misinformation on social media, and can also be used to combat other forms of online abuse and harassment, such as trolling and cyberbullying. Furthermore, given the significant amount of content on social media, the ability of these techniques to scale is limited. This is the point at which deep learning becomes useful. Deep learning is defined as a subset of machine learning used to obtain patterns and associations by processing and analyzing a large amount of data. In the context of online abuse and harassment, the use of deep learning can be used to automatically detect potentially offensive or hateful content, without the need for human interpretation or bias [5]–[11].

Although many popular SMPs started working in this direction and from fig. 1, we can observe how much hate content was removed by Facebook and Instagram since January 2020 to June 2022. Similarly, data statistics presented here for Twitter is from January 2020 to December 2021.

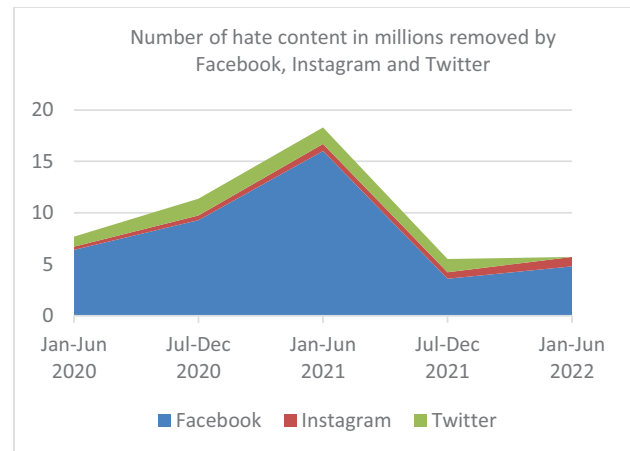


Fig. 1. Hate content removed by Facebook, Instagram, and Twitter Social Media Platforms

The data shown in the figure provides the intuition that how these social media communities are used as the vulnerable source for hate content dispersion. Considering only the six-month duration from July 2021 to December 2021, we observe that these three platforms removed approximately 5.53 million

contents. Besides the data removal policies adopted by SMPs, it has also been observed that these platforms have concentrated more on resourceful languages such as English, Spanish, etc. This is because the current research work for hate detection broadly supports popular languages such as English [2], [3], [12]–[14], Spanish [15]–[19], Turkish [6], [20], [21], etc. Given the limited good work available for the low-resource Hindi language, it became hard to detect hate comments on these platforms due to the support of multiple native languages for comment writing, which poses another hurdle in detection. Statistically, the Hindi language is the 4th primary spoken language globally, covering more than 330 million people. Furthermore, it is one of the two major official languages adopted by the Indian Constitution and covers around 46% population in India. The distribution of languages over the population of India is shown in fig. 2.

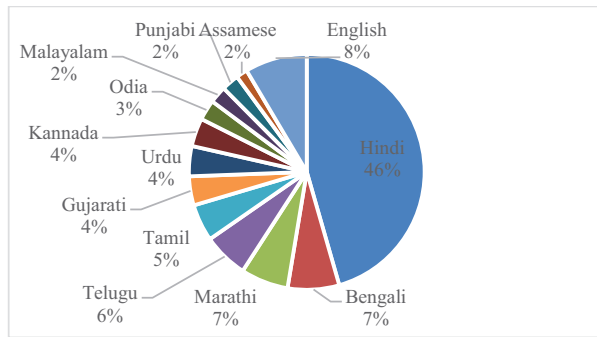


Fig. 2. Language versus speaker distribution in India

India is a land of diversity in terms of religion, living style, culture, language spoken, etc. Along with the various languages, each language has different dialects, which makes it more challenging to detect hate speech at ground level. Besides these wide varieties of languages, hate incidents are also increasing significantly every year. As per the data reported by the Government of India¹, the cases lodged under the hatred (Section 153A) are being increased yearly as shown in fig 3.

Although Hindi is a resource constraint language, its colossal user base serves as the salient source of information dissemination on social media. Fig. 3 shows how hate is disseminated using these platforms and makes society vulnerable to maintain peace and harmony.

This work focuses on automatic hate speech detection for the Hindi language. Often, hate speech is used to incite violence against a particular group. This type of speech is designed to stir up trouble and ultimately does nothing to improve the lives of those targeted.

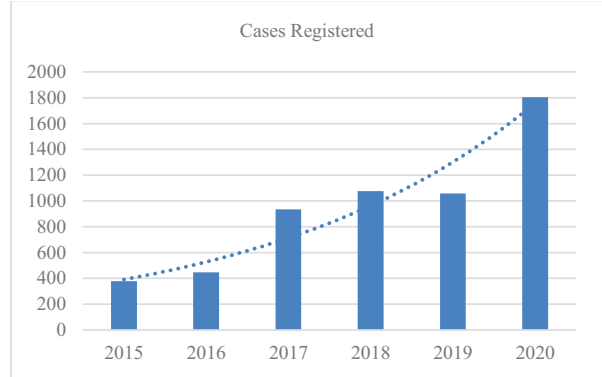


Fig. 3. Case registered under hate speech category in India.

It is meant to draw attention to oneself instead of cultivating positive change. However, speech intended to be hurtful sometimes results in violence against an individual or group. There is also various categorization of hate speech depending on the target, including sexual orientation, religious beliefs, race, and disability. An example of religious hate speech is shown in fig. 4.

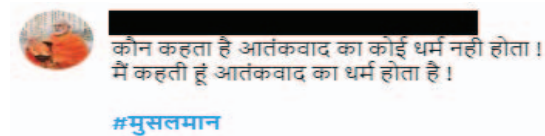


Fig. 4. Example of Hate speech in the Hindi language

This work proposes a two stage architecture using BERT encoder and CNN. The proposed model intends to detect hatred posts posted in Hindi language that threaten the dignity of individuals or communities. In this work, Hasoc 2020 [22] and Hasoc 2021 [23] Hindi datasets have been used to train and test the proposed model to ensure the generalizability of the model. The pre-trained BERT (Bidirectional Encoder Representations from Transformers) encoder model has been utilized to generate the embedding's from the textual data and to capture textual features; furthermore these embeddings are supplied as input to convolution neural network. A dropout layer is used to provide the learning stability to the model, and in the last, sigmoid layer has been used to generate the decision of hate or non-hate. Major contributions to this work are:

- Presented binary classification model for Hindi language to detect hate comments using BERT encoder and convolution neural network.
- Data augmentation using oversampling technique to equate majority and minority class.
- The result of proposed model has been evaluated on Hasoc 2020 & 2021 Hindi task dataset and compared

¹ <https://ncrb.gov.in/>

against transformer-based BERT classifier model for both balanced and imbalanced version of datasets.

The experiment's findings show that the suggested model performs better than the baseline model. The following sections of the paper are organized as follows: The relevant researches are covered in Section II; Section III presents the suggested model architecture. Section IV presents the experiments and their outcomes, while Section V presents the final remarks.

II. RELATED STUDIES

Social media has been a major player in information sharing in today's environment. It supports a better lifestyle to society, giving them a platform to share their knowledge, thoughts, experiences, etc. to benefit society, with better reach in short time span. Besides this, it also gives the power to miscreants to spread false narratives and to spread hate that leads to major violence sometimes. So it's better to identify these hate speeches on time so that they can be removed before they propagate nuisance in the society. Various researchers have proposed different models to detect hate speech automatically. These approaches solve the problems using meta information, linguistic features, and lexical resources [24].

Recent studies reported works on hate speech detection for the languages especially English [2], [3], [12]–[14], Turkish [20], [21], [25]–[27], and Arabic [1], [5], [7]. These studies reported various traditional machine learning models including SVM [28], Decision Tree, Naïve Bayes, and so on; advanced learning architectures such as CNN, LSTM, BiLSTM and transfer learning based models [29]–[31] and achieved better results for the relevant languages.

Few recent literatures also reported the hate speech detection for different versions of Hindi language such as Hindi-English code mixed language, roman Hindi language, and Hindi language. Agarwal et al. [4] identified hate speech in code mixed Hindi- English dataset and identified that many users were using Hindi slurs as well as English slurs to spread online swear. They have identified that use of native language slurs have contributed to the rise in code-switching languages used over social media platforms [7]. Another code-mixed Hindi-English dataset was presented by Bohra et al. [28] consisting tweets that were annotated at word level. They presented that SVM worked well and achieved an overall accuracy of 71% on their dataset. Moreover, Elouali et al. [8] utilized multiple datasets along with the code-mixing Hindi dataset and applied a character level CNN to detect hate speech, and they observed that hyper-tuning could significantly improve the performance of classifiers. Similarly, authors of [29] developed multiple classifiers over six public datasets and demonstrated that LR, BERT and RoBERTa classifiers performed equivalent on the Hindi dataset, while RoBERTa demonstrated the highest f1-score on the English dataset.

Phung & Cloos [30] implemented three Artificial Neural Network classifiers based on the skip-gram for Hindi, Bengali, and Hindi-Bengali. They used transfer learning for hate speech detection from classifiers trained on Hindi language and tested on Bengali language. Researchers have suggested that transfer learning can be applied to low computational models with

minor changes on the embedding layers to successfully apply these models to low resource languages. Three standard datasets of Hindi-English code-switched speech were tested against modified cross entropy based multi-lingual BERT and MuRIL models by the authors of [31]; results of the experiment revealed that modified cross-entropy leads to a significant improvement in f1-scores. Despite the good amount of work on HSD for Hindi language, we observed that the available models were trained on code-mixed datasets. At the same time limited studies such as [9], [10], [32], [33] are documented for Hindi language solely. This work aims to present an automated detector for detecting hate speech online in the low-resource Hindi language. A BERT-CNN model is proposed which takes text description as input and identifies whether the hate is present in that text or not. The internal architecture of the proposed model is described in the following section.

III. BERT-CNN MODEL

Automatic identification of hate content is a multi-domain problem following NLP and machine learning. It relies on feature engineering from the given text data and training a machine learning model, just like any other natural language text classification task. Although models dealing with sequential data such as RNN, LSTM has given an edge to solve hate detection task effectively, still they fail to deal with large volume data issue due to resource and computation constraints. Hence, transfer learning based pretrained model are used to handle this situation and tuned depending on the problem specific domain with little effort. Therefore, in this work, we have used transfer learning based BERT model along with convolution neural network to enhance the classification performance. The following subsection provides the details of the dataset followed by the architectural implementation of our proposed model.

A. Dataset

To train the proposed model, Hasoc 2020 [1] & Hasoc 2021 [2] datasets are used. These datasets contain the posts extracted from YouTube and Twitter. First dataset contains collection of hate and offensive posts in Hindi, English, German, Malayalam, and Tamil languages. From this dataset, we have used Hindi subtask dataset for this work. Similarly, Hasoc 2021 dataset supports Indo-Aryan and English languages. We considered Indo-Aryan subtask dataset for our work. Detailed categorical distribution of the data points is shown in table I, where one category represents for Hindi and other for mixed data points containing both Hindi and English.

TABLE I. DATASET DISTRIBUTION IN DETAILS

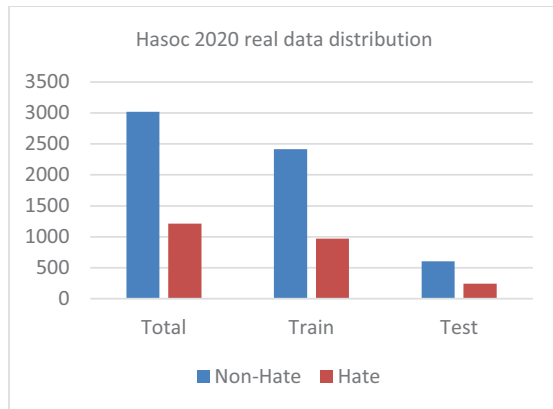
Dataset/ Type	<i>Hindi Tweets</i>	<i>Mixed Tweets</i>
Hasoc 2020	87.52%	12.48%
Hasoc 2021	86.48%	13.52%

The statistical distribution of classes, hate and non-hate is shown in fig 5. This distribution shows the original class distribution before oversampling followed by the oversampled version of the datasets.

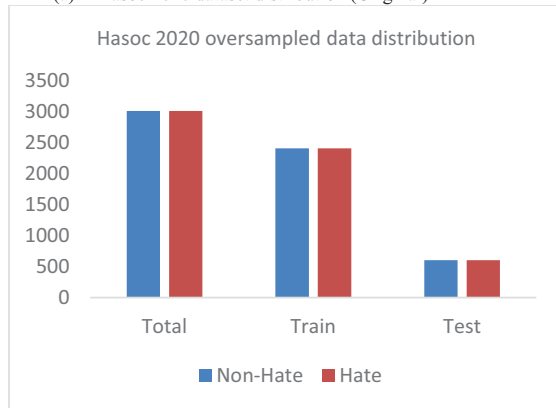
B. Preprocessing

The preprocessing is done as the initial step where cleaning of various unnecessary elements takes place. We retained the hashtags and emoji's as they convey meaningful information in the hate speech classification task as shown in fig. 4, where the hashtags tell us about the target community. In contrast, plain text post doesn't show any target information. The preprocessing step includes:

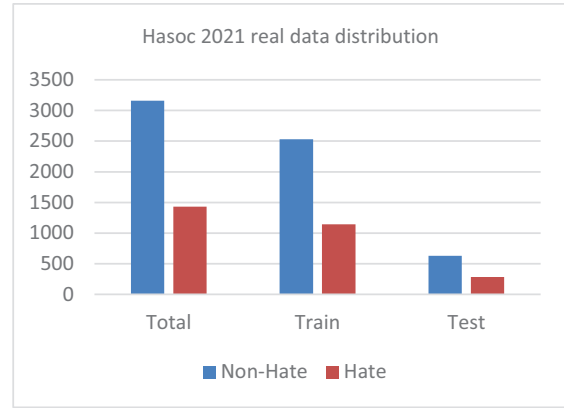
- URL's removal
- Usernames and mentions removal
- Punctuation removal
- Extra white space and long-period removal
- Translation of Emoticons into equivalent Hindi descriptions using a manually created translation dictionary.



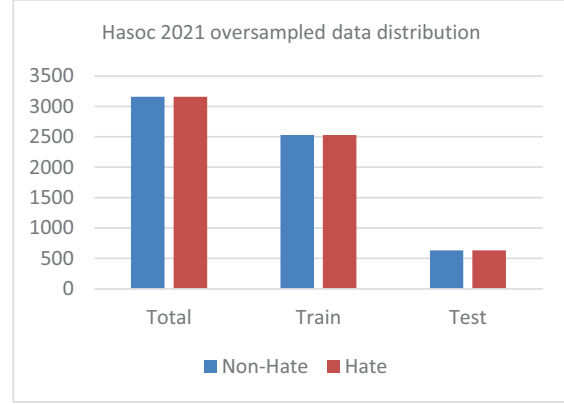
(a) Hasoc 2020 dataset distribution (Original)



(b) Hasoc 2020 dataset distribution (Oversampled)



(c) Hasoc 2021 dataset distribution (Original)



(d) Hasoc 2021 dataset distribution (Oversampled)

Fig. 5. Dataset Distribution

C. Pre-trained BERT Encoder

Pre-trained BERT Encoder is the first component of the proposed model, which takes the preprocessed output of the previous step and transforms this processed text data into contextualized embeddings. To obtain these contextual embeddings, we used the last four layers as our target output, which ultimately leads to an embedding vector of the size 768x4x64.

D. Convolution Neural Network

In next step, these embeddings are used in parallel to total of 128 convolution filters consisting of four different size convolution filters (786x1, 786x2, 786x3, 786x4) and 32 copies of each. Convolution filters are applied on the last four layers as individual channel operations, further convoluted outputs are passed through the ReLU activation function followed by the global pooling operation. In the end these pooled outputs concatenated together and passed to a flattened layer. This flattened output is further passed through a dense layer followed by the dropout and dense layer. The output of dense layer passed to sigmoid layer to generate final binary outcome. We used 20 epochs for training model. The architectural view of the proposed model is shown in fig 6.

IV. EXPERIMENTS AND RESULTS

The results of the proposed model were evaluated for the Hasoc 2020 & Hasoc 2021 dataset. From fig 5(a) & fig 5(c), we can observe that both datasets are imbalanced, so it will not be fair to use accuracy measure. Hence we evaluated our model performance using f1-score. It is defined as the harmonic mean of precision and recall; and gives us the tradeoff between high precision and high recall. Since both datasets are highly imbalanced, we performed oversampling for both datasets, and oversampled minority class to equate corresponding to majority class by creating copies of minority class. The distribution of oversampled dataset is shown in fig 5(b) & fig 5(d).

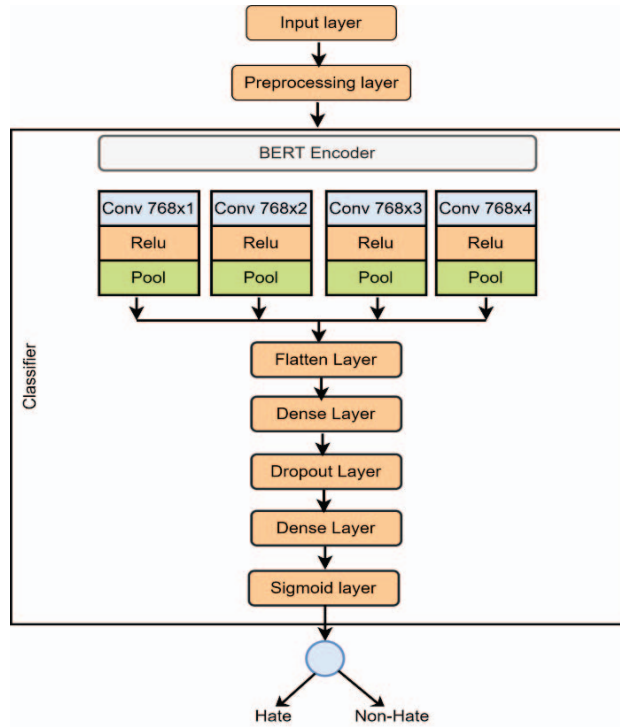


Fig. 6. Proposed BERT-CNN Model Architecture for hate speech detection in Hindi Language

Models used in this study were experimented with same dataset train-test split of 8:2 ratio. F1-Score metric was used for evaluation of these models. The results of experiments of investigated models are shown in Table II. The results show that the proposed work outperforms others with significant improvement in f1-score. The results have shown the great improvement on the oversampled dataset compared to imbalanced dataset as well as substantial improvement compared with the baseline BERT model.

TABLE II. F1-SCORE OF EXPERIMENTS PERFORMED ON HASOC 2020 AND HASOC 2021 DATASET

Model	Hasoc 2020		Hasoc 2021	
	<i>Imbalanced</i>	<i>Balanced</i>	<i>Imbalanced</i>	<i>Balanced</i>
BERT	0.49	0.80	0.52	0.72
BERT-CNN (ours)	0.51	0.84	0.54	0.77

V. CONCLUSION

Social media platforms have given the flexibility to the society to put their views without many regulatory conditions. This flexibility allows the hate speaker to be anonymous and hard to trace, which allows them to spread hate with less restrictions. Another central fact that made things complicated is the thin boundary present between hate speech and freedom of expression, which is not yet defined properly. In this work, oversampling method was used to handle the dataset imbalance. We proposed a combination of BERT encoder and Convolution neural network based BERT-CNN model to detect hate speech in Hindi language. To analyze the model performance, Hasoc 2020 and Hasoc 2021 Hindi datasets were used. The results were compared with baseline BERT model trained and tested on the same dataset. The model's performance as measured by the f1-score metric, demonstrates that our suggested model surpasses the base model with a notable f1-score improvement. The future direction for this work can be to allow other variants of BERT with different deep learning classification models to detect hate speech.

REFERENCES

- [1] M. Aljarah, Ibrahim and Habib, Maria and Hijazi, Neveen and Faris, Hossam and Qaddoura, Raneem and Hammo, Bassam and Abushariah, Mohammad and Alfawareh, "Intelligent detection of hate speech in Arabic social network: A machine learning approach," J. Inf. Sci., vol. 41, no. 4, pp. 483–501, 2021.
- [2] T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," Proc. 11th Int. Conf. Web Soc. Media, ICWSM 2017, no. Icwsm, pp. 512–515, 2017.
- [3] T. Davidson, D. Bhattacharya, and I. Weber, "Racial Bias in Hate Speech and Abusive Language Detection Datasets," pp. 25–35, 2019, doi: 10.18653/v1/w19-3504.
- [4] P. Agarwal, A. Sharma, J. Grover, M. Sikka, K. Rudra, and M. Choudhury, "I may talk in English but gaali toh Hindi mein hi denge : AA study of English-Hindi code-switching and swearing pattern on social networks," 2017 9th Int. Conf. Commun. Syst. Networks, COMSNETS 2017, no. January, pp. 554–557, 2017, doi: 10.1109/COMSNETS.2017.7945452.
- [5] H. Mohaouchane, A. Mourhir, and N. S. Nikolov, "Detecting Offensive Language on Arabic Social Media Using Deep Learning," 2019 6th Int. Conf. Soc. Networks Anal. Manag. Secur. SNAMS 2019, no. October, pp. 466–471, 2019, doi: 10.1109/SNAMS.2019.8931839.
- [6] C. Toraman and E. Halit Yilmaz, "Large-Scale Hate Speech Detection with Cross-Domain Transfer," arXiv Prepr. arXiv2203.01111, 2022.
- [7] O. Mohdeb, Djamilia and Laifa, Meriem and Zerargui, Fayssal and Benzaoui, "Evaluating transfer learning approach for detecting Arabic anti-refugee/migrant speech on social media," Aslib J. Inf. Manag., 2022.
- [8] A. Elouali, Z. Elberrichi, and N. Elouali, "Hate speech detection on multilingual twitter using convolutional neural networks," Rev. d'Intelligence Artif., vol. 34, no. 1, pp. 81–88, 2020, doi: 10.18280/ria.340111.
- [9] S. Kannan and J. Mitrović, "Hatespeech and Offensive Content Detection in Hindi Language using C-BiGRU Social Sentiment Analysis Financial IndeXes (ICT-15-2014 Grant: 645425) View project Europeana Libraries View project Hatespeech and Offensive Content Detection in Hindi Language using C-BiGRU," 2021.
- [10] R. Kumar and A. Kr Ojha, "KMI-Panlingua at HASOC 2019: SVM vs BERT for Hate Speech and Offensive Content Detection *," in CEUR Workshop Proceedings, 2021.
- [11] H. Faris, I. Aljarah, M. Habib, and P. A. Castillo, "Hate speech detection using word embedding and deep learning in the Arabic language context," ICPRAM 2020 - Proc. 9th Int. Conf. Pattern Recognit. Appl. Methods, no. March, pp. 453–460, 2020, doi: 10.5220/0008954004530460.

- [12] C. Science and P. Badjatiya, "Towards Identification , Classification and Analysis of Hate Speech on Social Media," no. June, 2019.
- [13] C. Harris, M. Halevy, A. Howard, A. Bruckman, and D. Yang, "Exploring the Role of Grammar and Word Choice in Bias Toward African American English (AAE) in Hate Speech Classification; Exploring the Role of Grammar and Word Choice in Bias Toward African American English (AAE) in Hate Speech Classification," in ACM Conference on Fairness, Accountability, and Transparency, 2022, vol. 22, pp. 789–798, doi: 10.1145/3531146.3533144.
- [14] S. Ghosh, Ekbal Asif, Pushpak Bhattacharyya, Saha Tisha, Kumar Alka, and S. Srivastava, "SEHC: A Benchmark Setup to Identify Online Hate Speech in English," IEEE Trans. Comput. Soc. Syst., 2022.
- [15] J. Haworth and P. Vincent, "Analysing Moral Beliefs for Detecting Hate Speech Spreaders on Twitter," in Advanced Geography and Geographical Learning, 2022, pp. 149–161.
- [16] A. Hasan, T. Sharma, A. Khan, M. Hasan, and A. Al-Abyadh, "Analysing Hate Speech against Migrants and Women through Tweets Using Ensembled Deep Learning Model," Comput. Intell. Neurosci., 2022, doi: 10.1155/2022/8153791.
- [17] P. Röttger, H. Seelawi, D. Nozza, Z. Talat, and B. Vidgen, "MULTILINGUAL HATECHECK: Functional Tests for Multilingual Hate Speech Detection Models," arXiv Prepr. arXiv2206.09917, 2022.
- [18] A. Arango, J. Pérez, B. Poblete, V. Proust, and M. Saldaña, "Multilingual Resources for Offensive Language Detection," in Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH), 2022, pp. 122–130.
- [19] J. Antonio García-Díaz, · Salud, M. Jiménez-Zafra, · Miguel, A. García-Cumbreras, and · Rafael Valencia-García, "Evaluating feature combination strategies for hate-speech detection in Spanish using linguistic features and transformers," Complex Intell. Syst., pp. 1–22, 2022, doi: 10.1007/s40747-022-00693-x.
- [20] Z. Melce, H. " Usünbeyius"usünbeyi, D. Akar, A. Arzucan", and A. A. Arzucan"ozgür, "Identifying Hate Speech using Neural Networks and Discourse Analysis Techniques," in Proceedings of the First Workshop on Language Technology and Resources for a Fair, Inclusive, and Safe Society within the 13th Language Resources and Evaluation Conference, 2022, pp. 32–41.
- [21] C. Toraman, H. Yilmaz, F. Şahinuç, and O. Ozcelik, "Impact of Tokenization on Language Models: An Analysis for Turkish; Impact of Tokenization on Language Models: An Analysis for Turkish," 2022.
- [22] T. Mandl et al., "Overview of the HASOC track at FIRE 2020: Hate speech and offensive content identification in Indo-European languages," CEUR Workshop Proc., vol. 2826, pp. 87–111, 2020.
- [23] T. Mandl et al., "Overview of the HASOC Subtrack at FIRE 2021: Hate Speech and Offensive Content Identification in English and Indo-Aryan Languages under Creative Commons License Attribution 4.0 International (CC BY 4.0)," 2021.
- [24] A. Schmidt and M. Wiegand, "A Survey on Hate Speech Detection using Natural Language Processing," Soc. 2017 - 5th Int. Work. Nat. Lang. Process. Soc. Media, Proc. Work. AFNLP SIG Soc., pp. 1–10, 2017, doi: 10.18653/V1/W17-1101.
- [25] C. Toraman, F. Şahinuç, and E. H. Yılmaz, "Large-Scale Hate Speech Detection with Cross-Domain Transfer," arXiv Prepr. arXiv2203.01111, 2022.
- [26] M. Y. PAK and S. GUNAL, "Impact of Text Representation and Preprocessing on Author Identification," ANADOLU Univ. J. Sci. Technol. A - Appl. Sci. Eng., vol. 18, no. 1, pp. 218–218, 2017, doi: 10.18038/aubtda.270276.
- [27] A. K. Uysal and S. Gunal, "The impact of preprocessing on text classification," Inf. Process. Manag., vol. 50, no. 1, pp. 104–112, 2014, doi: 10.1016/j.ipm.2013.08.006.
- [28] 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," Proc. 2nd Work. Comput. Model. PFopple's Opin. Personal. Emot. Soc. Media, PEOPLES 2018 2018 Conf. North Am. Chapter Assoc. Comput. Linguist. Hum. Lang. T, pp. 36–41, 2018, doi: 10.18653/v1/w18-1105.
- [29] N. Vashistha and A. Zubiaga, "Online multilingual hate speech detection: Experimenting with hindi and english social media," Inf., vol. 12, no. 1, pp. 1–16, 2021, doi: 10.3390/info12010005.
- [30] T. M. Phung and J. Cloos, "An exploratory experiment on Hindi, Bengali hate-speech detection and transfer learning using neural networks," no. NeurIPS, 2022.
- [31] A. Sharma, A. Kabra, and M. Jain, "Ceasing hate with MoH: Hate Speech Detection in Hindi-English code-switched language," Inf. Process. Manag., vol. 59, no. 1, p. 102760, 2022, doi: 10.1016/j.ipm.2021.102760.
- [32] V. K. Jha, P. Hrudya, N. V. Vinu, V. Vijayan, and P. Prabakaran, "DHOT-Repository and Classification of Offensive Tweets in the Hindi Language," Procedia Comput. Sci., vol. 171, pp. 2324–2333, 2020, doi: 10.1016/J.PROCS.2020.04.252.
- [33] A. Velankar, H. Patil, A. Gore, S. Salunke, and R. Joshi, "Hate and Offensive Speech Detection in Hindi and Marathi," in CEUR Workshop Proceedings, 2021.