

Analyzing Social Media Content for Detection of Offensive Text

Pawan Kalyan Jada¹, Konthala Yasaswini¹, Karthik Puranik¹,
Anbukkarasi Sampath², Sathiyaraj Thangasamy³ and Kingston Pal Thamburaj⁴

¹Indian Institute of Information Technology Tiruchirappalli

²Kongu Engineering College, Erode, Tamil Nadu, India

³Sri Krishna Adithya College of Arts and Science, Coimbatore

⁴Sultan Idris Education University, Tanjong Malim, Perak, Malaysia

Abstract

To tackle the conundrum of detecting offensive comments/posts which are considerably informal, unstructured, miswritten and code-mixed, we introduce two inventive methods in this research paper. Offensive comments/posts on the social media platforms, can affect an individual, a group or underage alike. In order to classify comments/posts in two popular Dravidian languages, Tamil and Malayalam, as a part of the HASOC - DravidianCodeMix FIRE 2021 shared task, we employ two Transformer-based prototypes which successfully stood in the top 8 for all the tasks. The codes for our approach can be viewed and utilized¹.

Keywords

Transformers, Sequence classification, Transliteration, Translation

1. Introduction

The term “Social media” provides a channel through which people engage in interactive communities and networks by creating, sharing, and exchanging thoughts and information. The growth rate of social media, especially Facebook and Twitter, has been exceptionally high since 2006. It has received users from almost all generations and all around the world. Users can interact and connect with others and form communities through social media. It allows users to share their ideas, views and information openly on various topics. This gives license to the users to write hateful and offensive comments sometimes. People come from a variety of racial backgrounds and hold a diversity of belief systems. This often causes for a conflict of opinions during their interactions on social media platforms. Many derogatory content target individuals based on their skin colour, gender, caste, nationality, religion, race, ethnicity. Due to the COVID-19 pandemic, the internet community has become more popular than it has ever been [1]. The amount of false narratives and derogatory remarks shared on online platforms has shot through the roof. A large number of social media users share malicious posts


¹<https://github.com/PawanKalyanJada/hasoc>

FIRE 2021: Forum for Information Retrieval Evaluation, December 13-17, 2021, India

✉ pawankj19c@iiitt.ac.in (P. K. Jada); konthalay18c@iiitt.ac.in (K. Yasaswini); karthikp18c@iiitt.ac.in (K. Puranik); anbu.1318@gmail.com (A. Sampath); sathiyarajt@skacas.ac.in (S. Thangasamy); fkingston@gmail.com (K. P. Thamburaj)



© 2021 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

despite understanding that they are infringing on their right to free expression. This may have a negative and adverse impact on user’s mental health. Various social media platforms restrict and minimize the profane comments by employing new rules and techniques.

Online hate speech and offensive content produces challenges to the society. The detection of hate speech is quite a daunting task, as the precise understanding of the speech largely depends on the circumstances it is being used in. Due to the extreme enormity of the internet and the growing number of online users, as well as the obscurity of the users, manually detecting and removing hate speech and profane information is a time-consuming and challenging job.

Code-mixing often entails the use of two languages to produce a third language that incorporates aspects from each in a functionally comprehensible manner. Low-resource languages such as Tamil and Malayalam are gaining significant attention alongside English on social networking platforms. The majority of the data on social media for these under-resourced languages is code-mixed. Tamil is a Dravidian language spoken by Tamils in India and Sri Lanka, as well as the Tamil community worldwide [2]. The official recognition of the language is in India, Sri Lanka, and Singapore. Tamil was the first to be classified as a classical language of India and is one of the longest-surviving classical languages in the world. Tamil has the oldest extant literature among Dravidian languages [3]. Malayalam is a Dravidian language spoken in southern India, having official language status in the Indian state of Kerala as well as the Union Territories of Lakshadweep and Puducherry. Malayalam scripts are alpha-syllabic, a type of “Abugida” writing system that is partially alphabetic and partly syllable-based.

This paper presents our work for the shared task on offensive language detection of code-mixed text in Dravidian languages (Malayalam-English and Tamil-English) at HASOC - DravidianCodeMix FIRE 2021. The rest of the paper is summarized as follows, 2 presents a discussion on the previous works on Offensive Language Detection in Dravidian Languages. 3 entails a detailed task description and analysis of the datasets for Tamil, Malayalam, and Kannada. In 4 we present a description of the models used for the tasks.

2. Related work

There has been a tremendous advancement in the research of offensive language detection over the past few years. On social media, hate speech in the form of racist and sexist statements is quite commonplace. In Waseem and Hovy, the authors provided a dataset of 16k tweets annotated for hate speech and analysed the features that help detect hate speech in the corpus. The authors of Davidson et al. used logistical regression to extract N-gram TF-IDF features from tweets and categorize each tweet into hate, offensive, and non-offensive categories. For identifying abusive language, the authors of Hassan et al. experimented with Support Vector Machines (SVMs) [7] trained on character and word-level features, Deep Neural Networks (DNNs) and Bidirectional Encoder Representations from Transformers (BERT) [8]. An Ensembling based approach which is based on hybridization of Naive Bayes, SVM, Linear Regression, and SGD classifiers was developed and tested on a Hindi-English code-mix dataset which outperformed the state-of-the-art systems and baseline models [9].

In Liu et al., the authors experimented with various classifiers which includes linear model with features of word unigrams, word2vec, and Hatebase; word-based Long Short-Term Memory

(LSTM) [11]; fine-tuned Bidirectional Encoder Representation from Transformer (BERT). Hande et al. created Kannada CodeMixed Dataset (KanCMD), a multitask learning dataset for sentiment analysis and offensive language identification. We work with several transformer-based models to classify social media comments as hope speech or not hope speech in English, Malayalam and Tamil languages. Various transformer-based models were fine-tuned to classify social media comments in English, Malayalam and Tamil languages into hope speech and non-hope speech labels [13, 14, 15]. In Yasaswini et al., the authors developed a model, CNN-BiLSTM, which has a layer of 1D convolutional layer followed by a dropout layer and then a bidirectional LSTM layer for identifying offensive language comments which are often code-mixed. The authors of Hande et al. introduced a Dual-Channel BERT4Hope approach employed by fine-tuning a language model based on BERT on the code-mixed data and its translation in English. Soft-voting is implemented on the fine-tuned transformer models to determine if any sentence contains information about an event that has occurred or not [18].

3. Task description and dataset

The main aim of the HASOC Shared task is to identify offensive content in the code-mixed comments/posts in the Dravidian languages collected from social media [19]. There are two tasks, in which task 1 is a message-level label classification task, systems have to classify a given YouTube comment in Tamil into offensive or not-offensive. Task 2 is also a message-level label classification task, in which systems have to classify a given tweet in code mixed Tamil and Malayalam into offensive or not-offensive.

We are provided with two different datasets for the two subtasks. The training dataset provided for task 1 comprised of 5877 YouTube comments in Tamil classified into offensive or not-offensive. The task 1 was limited to Tamil language. The dataset is also observed to be imbalanced. The training dataset provided for task 2 comprised of 4000 tweets in code mixed Tamil and 3999 tweets in code mixed Malayalam. The training and validation datasets provided for task 2 are well-balanced.

Language	Tamil (Task 1)	Tamil (Task 2)	Malayalam (Task 2)
not offensive	4,724	2,020	2047
offensive	1,153	1,980	1952
Total	5,877	4,000	3,999

Table 1

Class-wise distribution of the training set for both the Tasks

4. System Description

4.1. Task 1

We fine tune transformer based language models for this task. Firstly, emoticons and flags were cleaned from the dataset. Then the sentences were converted to lower case as some of the samples contains English text between them. We then pass the sequences through two

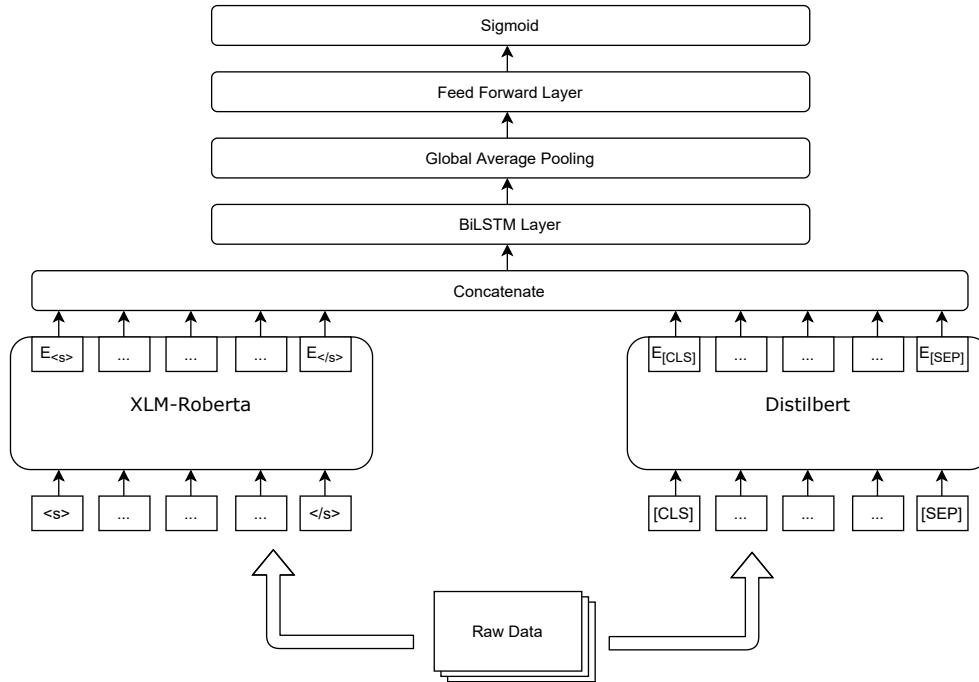


Figure 1: Model Architecture employed for Task 1

pre-trained models namely XLM-R and DistilBERT extracting the embeddings from both. These embeddings were then concatenated before being passed through the BiLSTM layers [20], eventually being pooled on a global average scale. These are fed to some Fully Connected layers and an Activation Function of sigmoid to get the probability scores as shown in Figure 1. By concatenating the embeddings, we expect that the model we created can benefit from knowledge of both the NLP models employed for the task, helping to distinguish better among the classes.

4.2. Task 2

For this task, we first transliterate [21] the Tamil sentences library in the English script to the native Tamil script by usage of “indic-transliteration” library¹. We then translate these sentences to English using the Google Translate API[22]. We then clean this parallel corpora of sentences by removing punctuations, stripping unwanted spaces at the end and converting the English sentences to a lower case. After the preprocessing, we tokenize these sequences using a tokenizer of a multilingual model, XLM-R. These tokens of Tamil and English are fed through the same XLM-R model and then passed through BiLSTM layers and a pooling layer at the end. Then we compute the weighted average of the Tamil and English vectors, with weights of 0.7 for Tamil and 0.3 for English. An Activation Function of sigmoid [23] is also applied at the end, deriving the probability scores required to classify a sentence. The entire architecture is shown

¹<https://pypi.org/project/indic-transliteration/>

in Figure 2. The same technique is done for Malayalam sentences as well, here the sentences being transliterated to Malayalam and the weights being 0.6 for Malayalam and 0.4 for English. These weights were set upon experimentation with various values and then selecting the best from all of them. Refer table 2 for parameters used in this task.

Parameters	Values
Number of LSTM layers	3
Number of LSTM units	128
Batch Size	32
Max Length	128
Optimizer	Adam
Learning Rate	1e-3
Activation Function	Sigmoid
Loss Function	cross-entropy

Table 2
Parameters used for training the model in Task 2

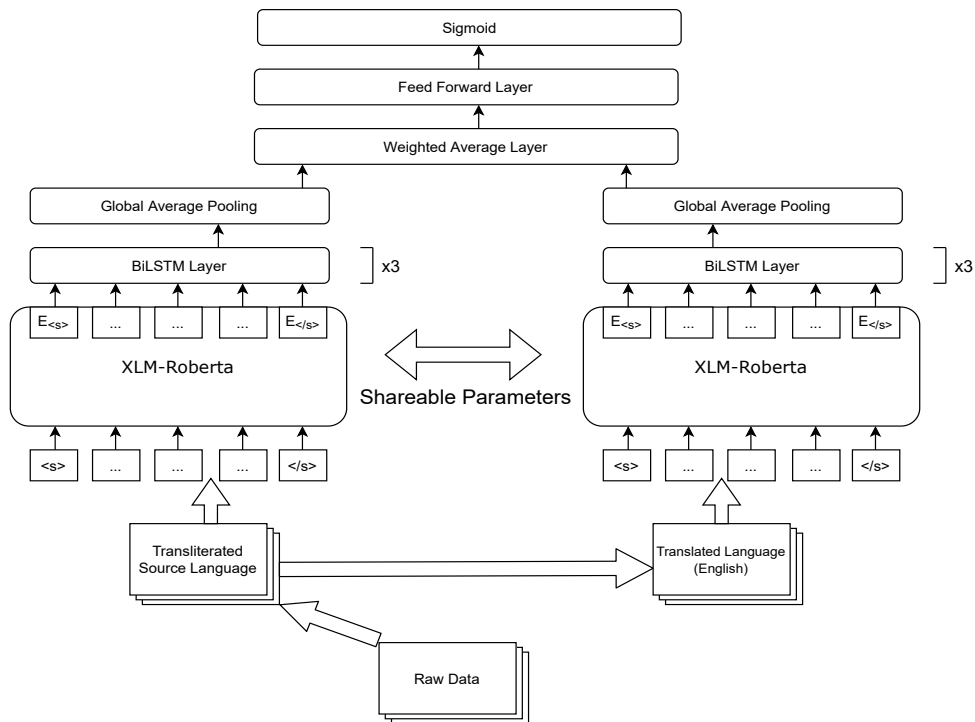


Figure 2: Model Architecture employed for Task 2

5. Methodology

5.1. XLM-RoBERTa

XLM-R [24] is a multilingual language model that achieved state-of-the-art results in all the multiple cross lingual benchmarks. One of the reason for the unparalleled performance is that it was trained on a mammoth 2.5 TB of CommonCrawl data [25]. It was trained with MLM loss as it's objective on 100 different languages, and it shares similar training routine as the one employed for RoBERTa [26] which is the reason the model is called XLM RoBERTa. XLM Roberta is fine-tuned for both of these tasks in a different architecture employed for the specific task. For this task we use *xlm-roberta-base* which consists of 12-layers, 768-hidden-state, 8-heads and a parameter size of 270M. The reason for selecting XLM-R is it outperformed other models in all the multilingual benchmarks.

5.2. DistilBERT

DistilBERT [27] is BERT's distilled version. It employs a triple loss language model that combines language modelling, distillation and cosine-distance losses. The two distillation losses in the triple loss have a significant influence on model performance. We fine tune *distilbert-base-multilingual-cased*, which is distilled from the mBERT checkpoint, for our cause in Task 1. It is known to have 40% less number of parameters than mBERT and runs 60% faster than it. The model consists of 6 layers, 768 dimensions, and 12 Attention heads, with a total of around 134 million parameters. We chose DistilBERT due to the less size of the model making it extract the word embeddings quicker.

6. Results

6.1. Task 1

The system description model on the Tamil language for this task gave a promising F1 score of 0.810. Embeddings from two of the most efficient multilingual pretrained models, XLM-R and DistilBERT were concatenated to extract their significant features. Further, the use of BiLSTM layers improves the accuracy as the information being fed doubles. The BiLSTM layers contains two LSTM's which take the input from the forward and backward directions and, thus, enhancing the context. However, one of the drawbacks which causes an impediment is the class imbalance between the "offensive" and "not-offensive" sentences in the test dataset. Also, XLM-R and DistilBERT follow BERT-based architectures and hence, the embeddings produced don't generate huge variations. Employment of pretrained models belonging to other architectures could produce higher accuracies. Increasing the number of models also might enhance the quality of embeddings produced, and thus, boosting the F1 scores further to 0.810.

6.2. Task 2

It was found out that our model gave an F1 score of 0.612 for Tamil and 0.670 for Malayalam. The dataset contains Tamil and Malayalam comments written in the Roman script, which is hard for

Task 1										
	Not offensive			Offensive			Overall			
Language	P	R	F	P	R	F	P	R	F	Acc
Tamil	0.888	0.875	0.882	0.468	0.500	0.484	0.812	0.807	0.810	0.807
Task 2										
	Not offensive			Offensive			Overall			
Language	P	R	F	P	R	F	P	R	F	Acc
Tamil	0.657	0.865	0.746	0.596	0.306	0.405	0.633	0.644	0.612	0.644
Malayalam	0.817	0.640	0.718	0.483	0.701	0.572	0.708	0.660	0.670	0.660

Table 3

Weighted F1 scores for task 1 and 2 by our system model on the test dataset where, P=Precision, R=Recall, F=F1-score, Acc=Accuracy

Task 1										
	Not offensive			Offensive			Overall			
Language	P	R	F	P	R	F	P	R	F	Acc
Tamil	0.868	0.891	0.879	0.500	0.446	0.471	0.796	0.804	0.799	0.804
Task 2										
	Not offensive			Offensive			Overall			
Language	P	R	F	P	R	F	P	R	F	Acc
Tamil	0.865	0.957	0.909	0.953	0.855	0.901	0.910	0.905	0.905	0.905
Malayalam	0.758	0.729	0.744	0.742	0.770	0.756	0.750	0.750	0.750	0.750

Table 4

Weighted F1 scores for task 1 and 2 by our system model on the validation dataset where, P=Precision, R=Recall, F=F1-score, Acc=Accuracy

the multilingual pretrained models trained on the native scripts to comprehend. Transliterating these sentences to the native language can prove to increase the F1 scores. Furthermore, we know that the models like XLM-R is trained on a large corpus of English sentences. Thus, the English translations of these transliterated dataset plays a huge role in further fine-tuning of the model. With the test dataset again containing sentences in native languages written in the Roman script, it was essential to give a higher precedence to the transliterated tokens over the translated tokens [28]. BiLSTM once again plays its role in ameliorating the results by increasing the information being fed.

However, we can never be definite of the accuracy in the transliterations and translations. Reduced quality if these sentences can affect the fine-tuning of the model significantly, and hence, lowering the F1 scores. Class imbalance prevails in this dataset too. With the ratio of not offensive to offensive in the range of 2:1, the model seems to find it arduous to predict the

offensive sentences efficiently. It is observed that the F1 scores between the not offensive and offensive differ by 0.15 to 0.3. This also impacts the overall F1 score for this task. However, as we can see in Table 4, the difference between the F1-scores of the “offensive” and “not-offensive” labels on the validation dataset didn’t seem to vary much. Table 3 tabulates the detailed weighted F1 scores for the test dataset.

7. Error Analysis

Few of the notable sentences where we felt that the sentences were misclassified have been discussed in this section. In task 1, 528 sentences were classified correctly in Tamil, while 126 failed to be classified well. There can be several reasons for this misclassification. Sometimes, the presence of offensive words doesn’t ensure that the sentence is “Offensive” and vice-versa. Comments are also filled with sarcasms, puns and typographical errors which have high probability of getting classified wrongly. Second task in Tamil has 645 sentences classified correctly and 356 wrongly. In Malayalam, 713 are correct and 238 wrong. We have discussed few of the Tamil sentences,

Task 1:

adey kirukka nalla paru,,,google unaku theriyuma,,, 2rs eppadi ellam pesura,,,Sanghis
This sentence is tagged as “not offensive”, but it is directed towards North Indians and probably as a reply to another comment/post.

Task 2:

Inta treylar kuta parkkira matiri illai.. Itai tiyettar la poy parkkanuma

The sentence is classified as offensive, but it is just a review which states, “this trailer itself isn’t good. Does someone have to go to the theatres too to watch this?”.

Another major drawback was the poor quality of translations by the Google API. The accuracy of classifications would have been better if the translations were of good quality in all the cases. For example,

Sentence 1: tl vere oru ss kandu

There is no other way is a very good translation of the Malayalam sentence and the model is able to learn and predict well.

Sentence 2: aga surya um jothikaum etho plan pani taga pola,not,Aga Surya Uma Jyotika is like something Plan Bani Daka

Aga Surya Um Jyotikaum Something like Plan snow is an example of how some sentences get partially got converted to English.

Sentence 3: aaiiii jolly yellam onnah polam onnah polam oannaa polam update app to view
IEE Jolly Yellam Onnah Bolam Onnah Bolam Oannaa Polam Uptade App To View is the supposed to be the English translation of the above sentence. We can see that the quality of

the translation is very bad.

8. Conclusion

Offensive language detection in social media posts presents to be a significant task due to social and marketing rationale. For the task of offensive language detection in code-mixed Dravidian languages Tamil and Malayalam, we introduce our research in this paper. For task 1, we extract the embeddings from XLM-R and DistilBERT, we concatenate them and pass them through BiLSTM layers. This model managed to give an F1 score of 0.810. Similarly, for task 2 we transliterate the dataset in which the Dravidian language is written in the Roman script and then, translate them into English and fine-tune the XLM-R model on it. This model gives us F1 scores of 0.612 for Tamil and 0.670 for Malayalam. Neglecting the fact that the English translations were of a poor quality, the model achieves very decent F1 scores for both the languages and, thus, opening a gateway for more research in this field.

References

- [1] A. Hande, K. Puranik, R. Priyadharshini, S. Thavareesan, B. R. Chakravarthi, Evaluating pretrained transformer-based models for covid-19 fake news detection, in: 2021 5th International Conference on Computing Methodologies and Communication (ICCMC), 2021, pp. 766–772. doi:10.1109/ICCMC51019.2021.9418446.
- [2] A. Hande, K. Puranik, R. Priyadharshini, B. R. Chakravarthi, Domain identification of scientific articles using transfer learning and ensembles, in: Trends and Applications in Knowledge Discovery and Data Mining: PAKDD 2021 Workshops, WSPA, MLMEIN, SD-PRA, DARAI, and AI4EPT, Delhi, India, May 11, 2021 Proceedings 25, Springer International Publishing, 2021, pp. 88–97.
- [3] B. R. Chakravarthi, V. Muralidaran, R. Priyadharshini, J. P. McCrae, Corpus creation for sentiment analysis in code-mixed Tamil-English text, in: Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL), European Language Resources association, Marseille, France, 2020, pp. 202–210. URL: <https://aclanthology.org/2020.sltu-1.28>.
- [4] Z. Waseem, D. Hovy, Hateful symbols or hateful people? predictive features for hate speech detection on Twitter, in: Proceedings of the NAACL Student Research Workshop, Association for Computational Linguistics, San Diego, California, 2016, pp. 88–93. URL: <https://aclanthology.org/N16-2013>. doi:10.18653/v1/N16-2013.
- [5] T. Davidson, D. Warmusley, M. Macy, I. Weber, Automated hate speech detection and the problem of offensive language, 2017. arXiv:1703.04009.
- [6] S. Hassan, Y. Samih, H. Mubarak, A. Abdelali, A. Rashed, S. A. Chowdhury, ALT submission for OSACT shared task on offensive language detection, in: Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection, European Language Resource Association, Marseille, France, 2020, pp. 61–65. URL: <https://aclanthology.org/2020.osact-1.9>.

- [7] M. Hearst, S. Dumais, E. Osuna, J. Platt, B. Scholkopf, Support vector machines, *IEEE Intelligent Systems and their Applications* 13 (1998) 18–28. doi:10.1109/5254.708428.
- [8] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. arXiv:1810.04805.
- [9] K. Yadav, A. Lamba, D. Gupta, A. Gupta, P. Karmakar, S. Saini, Bi-lstm and ensemble based bilingual sentiment analysis for a code-mixed hindi-english social media text, in: 2020 IEEE 17th India Council International Conference (INDICON), 2020, pp. 1–6. doi:10.1109/INDICON49873.2020.9342241.
- [10] P. Liu, W. Li, L. Zou, Nuli at semeval-2019 task 6: Transfer learning for offensive language detection using bidirectional transformers, in: Proceedings of the 13th international workshop on semantic evaluation, 2019, pp. 87–91.
- [11] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural computation* 9 (1997) 1735–80. doi:10.1162/neco.1997.9.8.1735.
- [12] A. Hande, R. Priyadharshini, B. R. Chakravarthi, KanCMD: Kannada CodeMixed dataset for sentiment analysis and offensive language detection, in: Proceedings of the Third Workshop on Computational Modeling of People’s Opinions, Personality, and Emotion’s in Social Media, Association for Computational Linguistics, Barcelona, Spain (Online), 2020, pp. 54–63. URL: <https://aclanthology.org/2020.peoples-1.6>.
- [13] K. Puranik, A. Hande, R. Priyadharshini, S. Thavareesan, B. R. Chakravarthi, Iiitt@l-edi-eacl2021-hope speech detection: There is always hope in transformers, 2021. arXiv:2104.09066.
- [14] B. R. Chakravarthi, HopeEDI: A multilingual hope speech detection dataset for equality, diversity, and inclusion, in: Proceedings of the Third Workshop on Computational Modeling of People’s Opinions, Personality, and Emotion’s in Social Media, Association for Computational Linguistics, Barcelona, Spain (Online), 2020, pp. 41–53. URL: <https://aclanthology.org/2020.peoples-1.5>.
- [15] B. R. Chakravarthi, V. Muralidaran, Findings of the shared task on hope speech detection for equality, diversity, and inclusion, in: Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion, Association for Computational Linguistics, Kyiv, 2021, pp. 61–72. URL: <https://aclanthology.org/2021.ltedi-1.8>.
- [16] K. Yasaswini, K. Puranik, A. Hande, R. Priyadharshini, S. Thavareesan, B. R. Chakravarthi, IIIT@DravidianLangTech-EACL2021: Transfer learning for offensive language detection in Dravidian languages, in: Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages, Association for Computational Linguistics, Kyiv, 2021, pp. 187–194. URL: <https://aclanthology.org/2021.dravidianlangtech-1.25>.
- [17] A. Hande, R. Priyadharshini, A. Sampath, K. P. Thamburaj, P. Chandran, B. R. Chakravarthi, Hope speech detection in under-resourced kannada language, 2021. arXiv:2108.04616.
- [18] P. Kalyan, D. Reddy, A. Hande, R. Priyadharshini, R. Sakuntharaj, B. R. Chakravarthi, IIIT at CASE 2021 task 1: Leveraging pretrained language models for multilingual protest detection, in: Proceedings of the 4th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2021), Association for Computational Linguistics, Online, 2021, pp. 98–104. URL: <https://aclanthology.org/2021.case-1.13>. doi:10.18653/v1/2021.case-1.13.
- [19] B. R. Chakravarthi, P. K. Kumaresan, R. Sakuntharaj, A. K. Madasamy, S. Thavareesan,

- P. B, S. Chinnadayar Navaneethakrishnan, J. P. McCrae, T. Mandl, Overview of the HASOC-DravidianCodeMix Shared Task on Offensive Language Detection in Tamil and Malayalam, in: Working Notes of FIRE 2021 - Forum for Information Retrieval Evaluation, CEUR, 2021.
- [20] G. Xu, Y. Meng, X. Qiu, Z. Yu, X. Wu, Sentiment analysis of comment texts based on bilstm, *Ieee Access* 7 (2019) 51522–51532.
- [21] K. Regmi, J. Naidoo, P. Pilkington, Understanding the processes of translation and transliteration in qualitative research, *International Journal of Qualitative Methods* 9 (2010) 16–26.
- [22] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, Łukasz Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, J. Dean, Google’s neural machine translation system: Bridging the gap between human and machine translation, 2016. [arXiv:1609.08144](https://arxiv.org/abs/1609.08144).
- [23] X. Yin, J. Goudriaan, E. A. Lantinga, J. Vos, H. J. Spiertz, A flexible sigmoid function of determinate growth, *Annals of botany* 91 (2003) 361–371.
- [24] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, V. Stoyanov, Unsupervised cross-lingual representation learning at scale, *arXiv preprint arXiv:1911.02116* (2019).
- [25] J. Smith, H. Saint-Amand, M. Plamadă, P. Koehn, C. Callison-Burch, A. Lopez, Dirt cheap web-scale parallel text from the common crawl, in: *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2013, pp. 1374–1383.
- [26] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert pretraining approach, *arXiv preprint arXiv:1907.11692* (2019).
- [27] V. Sanh, L. Debut, J. Chaumond, T. Wolf, Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, *arXiv preprint arXiv:1910.01108* (2019).
- [28] K. Puranik, A. Hande, R. Priyadharshini, T. Durairaj, A. Sampath, K. Thamburaj, B. R. Chakravarthi, Attentive fine-tuning of transformers for translation of low-resourced languages @loresmt 2021, 2021.