

Multilingual Models for Sentiment and Abusive Language Detection for Dravidian Languages

Anand Kumar M

Department of Information Technology
National Institute of Technology Karnataka, Surathkal
India
m_anandkumar@nitk.edu.in

Abstract

This work delves into the realm of abusive comment detection and sentiment analysis within code-mixed content, focusing specifically on Dravidian languages. The languages covered include Tulu, and Tamil. For this investigation, TFIDF-based Long Short-Term Memory (LSTM) and Hierarchical Attention Networks (HAN) are employed as the analytical tools.

Interestingly, the research highlights the prevalence of traditional TF-IDF techniques over Hierarchical Attention models in both sentiment analysis and the identification of abusive language across the diverse linguistic landscape encompassing Tulu and Tamil.

Of note is the Tulu sentiment analysis system, which demonstrates remarkable prowess in handling Positive and Neutral sentiments. In contrast, the sentiment analysis system tailored for Tamil exhibits comparatively lower performance levels. This discrepancy underscores the critical need for well-balanced datasets and intensified research endeavors to enhance the accuracy of sentiment analysis, particularly in the context of the Tamil language.

Shifting focus to abusive language detection, the TF-IDF-LSTM models consistently outperform the Hierarchical Attention models. Intriguingly, the mixed models exhibit particular strength in classifying categories like "Homophobia" and "Xenophobia." This intriguing outcome accentuates the value of incorporating both code-mixed and original script data, presenting novel avenues for advancing social media analysis research in diverse linguistic scenarios involving the Dravidian languages.

1 Introduction

The number of users is exponentially increasing on online social media platforms daily. More than 4.74 billion people used social media platforms¹ in the year 2022. Furthermore, the number of users

¹<https://influencermarketinghub.com/social-media-sites/>

will continue to grow even higher with cheaper internet and smartphones. Many online abusers use social media platforms as a venue to abuse other users through comments or posts. Nowadays, the marketing industry heavily relies on social media comments posted by users about their products. On the other hand, political parties base their political movements on the opinions expressed by citizens on social media. Government policies are revised based on the sentiments of the citizens identified through social media. Therefore, analyzing the comments posted on social media platforms is the most trending research domain in Natural Language Processing. These social media comments have opened up new and exciting research directions for NLP.

In a multilingual country like India, mixing languages while speaking is a typical behavior of the people. However, many people do not mix languages while writing for general purposes. However, this trend has changed in the era of social media, and people tend to mix languages when posting comments on online platforms. Users mainly use the Roman script to write comments, even in their native language. This phenomenon is known as code-mixing.

This paper presents a system developed for abusive language detection and sentiment analysis tasks conducted at the DravidianLangtech-2023. We have developed three different systems to identify abusive text and sentiment in social media posts. The methods used are the Hierarchical attention-based LSTM, TFIDF-based LSTM, and mixed language model. Additionally, to address the data imbalance, we have used contextualized embedding-based text generation to generate comments for the minority class.

2 Related Works

Recently, there has been a considerable amount of work and effort to collect resources for code-

switched text in various languages. However, code-switched datasets and lexicons specifically for sentiment analysis purposes are still limited in number, size, and availability (Chakravarthi et al., 2018, 2019a,b,c; Padmamala and Vijayarani, 2017; Ranjan et al., 2016; Ar et al., 2012; Devi and Kanimuthu, 2023).

For monolingual sentiment analysis, various corpora have been developed for different languages. For example, the work by (Wiebe et al., 2005) introduced an annotated corpus for sentiment analysis in English. Similarly, the Rusentiment corpus was created for sentiment analysis in Russian (Rogers et al., 2018), the Twitter Sentiment Corpus (TSC) was developed for sentiment analysis in German (Cieliebak and Diab, 2017), and the Norwegian Social Media Corpus (NoReC) was annotated for sentiment analysis in Norwegian (Mæhlum et al., 2019).

In the context of code-mixing, several datasets have been created to facilitate sentiment analysis. (Sitaram et al., 2015; Joshi et al., 2016; Patra et al., 2018) worked on building an English-Hindi corpus for sentiment analysis. (Solorio et al., 2014; Vilares et al., 2015, 2016) introduced an English-Spanish corpus for sentiment analysis. (Lee et al., 2015) collected a Chinese-English corpus from Weibo.com for sentiment analysis, and (Patra et al., 2018) released English-Bengali data for sentiment analysis.

Tamil, a Dravidian language spoken by Tamil people in India, Sri Lanka, and the Tamil diaspora, has received attention in sentiment analysis research (Padmamala and Vijayarani, 2017). The growing number of native Tamil speakers presents a potential market for commercial NLP applications (Ranjan et al., 2016). However, sentiment analysis on Tamil-English code-mixed data is relatively underdeveloped, and readily available data for research purposes is limited.

In the past, research on code-mixed corpora primarily relied on word-level annotations. However, this approach is not only time-consuming but also expensive to create. To address this limitation, researchers have explored the use of neural networks and meta-embeddings, which have shown promise in code-switched research without the need for word-level annotation (Kiela et al., 2018; Winata et al., 2019c).

(Winata et al., 2019a) demonstrated the effectiveness of utilizing information from pre-trained embeddings without explicit word-level language

tags in code-switched sentiment analysis. This approach leverages the power of neural networks to learn representations that can capture sentiment in code-mixed data.

Furthermore, (Winata et al., 2019b) introduced a method to utilize subword-level information from closely related languages to enhance the performance of sentiment analysis on code-mixed text. By leveraging the linguistic similarities between languages, this approach aims to improve the accuracy of sentiment analysis in code-mixed data.

In this field, there has not been extensive Tamil language-oriented research. One important reason for this could be the scarcity of data in social media in Tamil compared to other languages, especially English, and the limited availability of linguistic resources in Tamil. Many datasets have been created in Tamil to promote more research in this language. One of them is "HopeEDI" (Equality, Diversity, and Inclusion), a dataset for hope speech in Tamil (Chakravarthi et al., 2020). Several baselines have also been created to standardize the dataset.

Research on abusive comment detection in Tamil is still in its early stages, but it has made significant progress in recent years. The earliest models on text classification used linear classifiers. This was followed by several works based on Deep Learning methods. Recurrent Neural Networks like LSTMs showed promising results. (Mandalam and Sharma, 2021) classified Dravidian Tamil and Malayalam code-mixed comments according to their polarity and used the LSTM architecture. In (Arora, 2020), a pre-trained version of ULM-FiT was used to develop a model to detect hate speech in Tamil-English social media comments.

This was followed by the use of transformer-based models after being introduced in (Vaswani et al., 2017), and further exploration was done after the release of BERT (Devlin et al., 2019). In (Mishra and Mishra, 2019), MultiLingual BERT and monolingual BERT were used for hate speech identification in Indo-European languages. In (Ziehe et al., 2021), the authors fine-tuned XLM-RoBERTa (Conneau et al., 2020) for Hope Speech Detection in English, Malayalam, and Tamil texts. Recently, in (García-Díaz et al., 2022), the authors proposed a method for detecting abusive comments in Tamil using multilingual transformer models. And in (Prasanth et al., 2022), they performed abuse detection using TF-IDF and the Random Kitchen Sink Algorithm on Tamil text.

3 Dataset Description

We utilized a dataset for sentiment analysis and detection of abusive language, which was sourced from the (Priyadharshini et al., 2023) and (Hegde et al., 2023) references. This dataset was employed as part of the shared task held during the third Dravidian Lang Tech workshop at RANLP-2023. The dataset provided by the organizers encompassed content in Code-Mixed Tamil, as well as Tamil and Telugu languages. The data was curated from various social media interactions, such as posts and comments.

In terms of annotation, the Code-Mixed Tamil and Tamil comments were assigned labels from a set of 8 categories: None, Misandry, Misogyny, Xenophobia, Homophobia, Transphobia, Hope Speech, and Counter Speech. Conversely, the Telugu comments were categorized into just two classes: Hate and non-Hate. The data samples with specific labels can be found in Tables 1 and 2 for your reference.

Notably, the Code-Mixed Tamil dataset contained a larger number of comments compared to the other two datasets. For evaluation purposes, approximately 20% of the data was allocated for testing in both the Tamil datasets. Within the datasets, nearly 50% of the content fell under the "None" class. It's important to highlight that the imbalanced distribution was due to the greater number of classes present in the Code-Mixed Tamil and Tamil comments. On the contrary, the Telugu dataset exhibited a balanced distribution, and no validation data was included for this dataset. The distribution of data for training, validation, and testing is presented in detail in Table 3. The prevalent class across posts was "None," followed by the "Misandry" class within the Tamil dataset. In the context of Telugu, there were 1939 posts labeled as Hate and 2061 as Non-Hate.

For the sentiment analysis task, the organizers furnished social media comments in both Tamil and Tulu languages. The Tamil dataset was annotated with four sentiment classes: positive, negative, mixed, and unknown. Similarly, the Tulu dataset encompassed four classes: positive, negative, neutral, and unknown. The training set for Tamil sentiment analysis contained around 34,000 comments, while the Tulu dataset comprised 6674 posts. In the Tamil training dataset, the positive class accounted for 20,000 posts, with the remaining classes containing between 4,000 to 5,000 posts.

In the Tulu dataset, around 3,000 posts were labeled as positive and 1,800 as neutral. As a noteworthy point, since the Tamil dataset featured an unknown class and the Tulu dataset contained a neutral class, these two classes were considered equivalent within the context of language mixed models.

4 Methodology

4.1 TFIDF-LSTM

We utilized traditional and robust TF-IDF models to generate term vectors. These term vectors serve as embeddings for each word, for example the term "loose" in the example "Loose kooda interveiw panreenga kuruttu koothikku innoru .." would be represented differently based on the context for a Tamil it would be matched with the vector for crazy while in an English sentence it would retain its original vector, and TF-IDF vectors are created for each post. The learned TF-IDF vectors for each post were then inputted into a BiLSTM (Bidirectional LSTM) to further capture contextual information in both directions. The BiLSTM layer was composed of 100 units, transforming the context vector for each post. Finally, we employed a machine learning classifier to classify unseen posts. During the validation process, we discovered that the Linear SVM model provided the best results compared to other models. We employed the TweetTokenizer to tokenize the code-mixed posts and jointly learned the character and word n-gram models up to the trigram level to acquire the vectors. We determined the optimal parameters for the model using grid search.

4.2 Hierarchical Attention Networks

We experimented with Hierarchical attention-based LSTM models to capture the latent information from the code-mixed comments. Since the dataset is derived from social media sources, we incorporated character-level embeddings in the first layer of the Hierarchical attention network. By using the character sequence, we learned individual word vectors. This approach entails initially learning the words from the characters, followed by combining the word vectors to form the embedding for each post or comment. We employed attention mechanisms to assign importance to specific words within the post, for example in the sentence "Loose kooda **interveiw panreenga kuruttu koothikku** innoru ..", the terms that focuses on the gender and

Table 1: Language, Task, Examples and its corresponding Labels

Language	Task	Examples	Labels
Tamil	Sentiment	Ithu yethu maathiri illama puthu maathiyaala irukku	Positive
		Waste padam tharu maru flop aamai nakkis	Negative
Tulu	Sentiment	Irena tulu ucharane bhari likundu	Positive
		Ayana pukuli n ora nadt korle...	Negative
Telugu	Hate Speech	Torch lite Kuda Leni rojula fake news vadu kavalane chesind	hate
		Mallareddy Dookudu cenima lo bramhi character correct set aithadu	non-hate
Tamil Code-mixed	Misogyny	poda nee oruru punda yechakala raja	Misandry
		Loose kooda interveiw panreenga kuruttu koothikku innoru ..	Misogyny
		Entha Mari aravaningala seruppala adikkanum entha Mari prachanai...	Transphobic

Table 2: Language, Task and its corresponding Labels

Language	Task	Labels
Code-Mixed Tamil	Abusive Lang. Detection	None, Misandry, Misogyny, Xenophobia,
Tamil	Abusive Lang. Detection	Homophobia, Transphobia, Hate Speech and Counter
Telugu	Abusive Lang. Detection	Hate and non-Hate
Tamil	Sentiment Analysis	Positive, Negative, Mixed and Unknown
Tulu	Sentiment Analysis	Positive, Negative, Neutral and Unknown

Table 3: Dataset distribution for Train, Test, and Validation sets

Language	Task	Train Set	Validation set	Test set
Code-Mixed Tamil	Abusive Lang. Detection	5948	1488	1857
Tamil	Abusive Lang. Detection	2240	560	699
Telugu	Abusive Lang. Detection	4000	-	500
Tamil	Sentiment Analysis	33990	3787	650
Tulu	Sentiment Analysis	6674	903	749

actions which relate to "Misogyny", have more attention weights than the other terms. In the case of abusive language and sentiment detection, certain words in social media comments have a significant impact on determining the type of abuse and sentiment. Hence, we utilized the hierarchical attention network to capture the underlying information. Additionally, we employed Bi-LSTM for learning the sequence vectors.

4.3 Multilingual Models

In the sentiment analysis model, we combined the Tulu and Tamil code-mixed datasets since both languages belong to the Dravidian language family and share common English words in their code-mixed posts. We trained a multilingual sentiment

analysis model using the previously proposed TF-IDF-based Bi-LSTM models on the mixed language dataset. This model was trained once and then tested separately for Tulu and Tamil sentiment analysis. We hypothesized that the shared features between the two languages could assist each other in the sentiment analysis task.

In our pursuit of abusive language detection, we took an innovative step by merging the Tamil code-mixed dataset with the authentic Tamil dataset. The purpose behind this combination was to train a single model that could effectively identify abusive language. This approach aimed to explore how the amalgamation of code-mixed and pure script data could influence the model's performance in detecting abusive language. Additionally, English swear

words are often mixed with regional social media posts. Users may post content in their regional language but incorporate English swear words to abuse someone. We believed that the mixture of both datasets would provide a unique research perspective in social media analysis. Importantly, in the future, such mixed models will become increasingly important, as opposed to relying solely on language-specific tools.

Although the Telugu abusive dataset was provided, we did not combine it with the other datasets due to the mismatch in classes. The Telugu dataset consists of only two classes: hate and not hate. These mixed language learning approaches draw inspiration from code-mixed transfer learning-based POS tagging methods (Madasamy and Padannayil, 2021).

5 Results and analysis

In this section, we discuss the results obtained for the proposed models, as shown in Tables 4 and 5. Generally, the Hierarchical Attention models did not outperform the traditional TF-IDF-based techniques.

In the sentiment analysis task, the accuracy and F1 score for the Tamil dataset were relatively low for all the developed methods. This could be due to the highly imbalanced nature of the dataset. Upon analyzing the class-specific performance of the Tamil sentiment analysis, we found that the recall was higher for the positive class, while the precision was higher for the negative class compared to the other precision and recall values. Although the mixed language models did not perform significantly better, they exhibited higher precision for the negative class and higher recall for the positive class compared to the TF-IDF-LSTM model. A similar trend was observed in the Tulu sentiment analysis system. Additionally, the Tulu models performed better for the Positive and Neutral classes, with an F1 score of 0.83 for positive class and 0.63 for neutral class. The macro F1 score for the mixed language model was 0.47, whereas for the TF-IDF-LSTM model, it was 0.52.

In the abusive language detection task for Telugu, the TF-IDF-LSTM models outperformed the Hierarchical Attention models. The F1 score for the hate class was 0.62 and for the non-hate class it was 0.66. For the Tamil test set, the macro F1 score for the TF-IDF-LSTM model was the same as that of the mixed model. When analyzing the class-wise

performance, the mixed models performed better in the "Homophobia" and "Xenophobia" classes. The HAN model failed to detect certain classes with fewer training posts, indicating that the HAN model requires more comments to train effectively. For the Tamil code-mixed abusive language detection task, the mixed models performed better for the "Counter Speech" class and "Misandry". Overall, the recall of the mixed models was comparable to the TF-IDF model, but they exhibited lower precision.

6 Conclusion and Future Scope

The traditional TF-IDF-based techniques have outperformed the Hierarchical Attention models in both the sentiment analysis and abusive language detection tasks. This suggests that, for the given datasets and tasks, the TF-IDF approach provided superior results.

The sentiment analysis task for the Tamil dataset exhibited lower accuracy and F1 scores, which may be attributed to the highly imbalanced nature of the dataset. When examining the class-specific performance, it was found that the positive class had higher recall while the negative class had higher precision. This indicates the need for addressing the dataset imbalance to improve the overall performance of sentiment analysis models.

Although the mixed language models did not show significant improvements, they displayed some advantages compared to the TF-IDF-LSTM model. These models exhibited higher precision for the negative class and higher recall for the positive class in both sentiment analysis and abusive language detection tasks. This suggests that leveraging mixed language data could be beneficial, and further exploration and enhancement of these models are warranted.

In conclusion, this work provides insights into the challenges and potential improvements in sentiment analysis and abusive language detection for code-mixed data, specifically focusing on Dravidian languages. Future work should address dataset imbalance, explore enhanced mixed language models, expand datasets, improve models through advanced architectures, adopt a multilingual approach, and investigate fine-tuning and transfer learning techniques. By tackling these areas, researchers can enhance the performance and robustness of sentiment analysis and abusive language detection in code-mixed scenarios, con-

Table 4: Results for Tamil and Tulu Language Models

Language	Models	P	R	F1	Accuracy
Tamil	TFIDF-LSTM	0.31	0.29	0.24	0.25
	HNN	0.20	0.20	0.20	0.20
	Mixed	0.29	0.25	0.18	0.18
Tulu	TFIDF-LSTM	0.55	0.51	0.51	0.67
	HNN	0.34	0.23	0.23	0.54
	Mixed	0.49	0.47	0.47	0.63

Table 5: Results for Tamil and Tulu Code-Mixed Language Models

Language	Models	P	R	F1	Acc
Telugu	TFIDF-LSTM	0.65	0.64	0.64	0.64
	HNN	0.49	0.49	0.49	0.49
	Mixed	-	-	-	-
Tamil	TFIDF-LSTM	0.44	0.33	0.35	0.69
	HNN	0.23	0.23	0.23	0.64
	Mixed	0.43	0.32	0.35	0.67
Code-Mixed Tamil	TFIDF-LSTM	0.64	0.45	0.51	0.74
	HNN	0.34	0.23	0.23	0.70
	Mixed	0.60	0.44	0.49	0.73

tributing to the advancement of natural language processing in diverse linguistic contexts.

Acknowledgements

We thank the Sentiment Analysis and Abusive language Detection Shared Task Organizers.

References

- Ortal Ar, Moshe Koppel, Dirk Börner, and Dana Soffa. 2012. Cross-lingual sentiment analysis for languages with scarce resources. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*.
- Gaurav Arora. 2020. [Gauravarora@hasoc-dravidian-codemix-fire2020: Pre-training ulmfit on synthetically generated code-mixed data for hate speech detection](#).
- B. R. Chakravarthi, M. Arcan, and J. P. McCrae. 2018. Improving wordnets for under-resourced languages using machine translation. In *Proceedings of the 9th Global WordNet Conference (GWC 2018)*, page 78.
- B. R. Chakravarthi, M. Arcan, and J. P. McCrae. 2019a. Comparison of different orthographies for machine translation of under-resourced dravidian languages. In *2nd Conference on Language, Data and Knowledge (LDK 2019)*.
- B. R. Chakravarthi, M. Arcan, and J. P. McCrae. 2019b. Wordnet gloss translation for under-resourced languages using multilingual neural machine translation. In *Proceedings of the Second Workshop on Multilingualism at the Intersection of Knowledge Bases and Machine Translation*, pages 1–7, Dublin, Ireland. European Association for Machine Translation.
- B. R. Chakravarthi, R. Priyadharshini, B. Stearns, A. S. Jayapal, M. Arcan, M. Zarrouk, and J. P. McCrae. 2019c. Multilingual multimodal machine translation for dravidian languages utilizing phonetic transcription. In *Proceedings of the 2nd Workshop on Technologies for MT of Low Resource Languages*, pages 56–63, Dublin, Ireland. European Association for Machine Translation.
- Bharathi Raja Chakravarthi, Vigneshwaran Muralidaran, Ruba Priyadharshini, and John Philip McCrae. 2020. [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)*, pages 202–210, Marseille, France. European Language Resources association.
- M. Cieliebak and M. Diab. 2017. Twitter language identification of arabic-english code-switching. In *Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching*, pages 1–11, Vancouver, Canada. Association for Computational Linguistics.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#).

- Sharmila Devi and S. Kannimuthu. 2023. [Author profiling in code-mixed whatsapp messages using stacked convolution networks and contextualized embedding based text augmentation](#). *Neural Process. Lett.*, 55(1):589–614.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- José García-Díaz, Manuel Valencia-García, and Rafael Valencia-García. 2022. [UMUTeam@TamilNLP-ACL2022: Abusive detection in Tamil using linguistic features and transformers](#). In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 45–50, Dublin, Ireland. Association for Computational Linguistics.
- Asha Hegde, Bharathi Raja Chakravarthi, Rahul Shashirekha, Hosahalli Lakshmaiah and Ponnusamy, SUBALALITHA CN, Lavanya S K, Thenmozhi D, Shreya Karunakar, Martha and Shreeram, and Sarah Aymen. 2023. Findings of the shared task on sentiment analysis in tamil and tulu code-mixed text. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.
- Aditya Joshi, Monojit Choudhury, and Mark J Carman. 2016. Towards building a code-mixed social media corpus for three indian languages. In *Proceedings of the 10th International Conference on Language Resources and Evaluation (LREC)*.
- D. Kiela, C. Wang, and K. Cho. 2018. Dynamic meta-embeddings for improved sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1466–1477, Brussels, Belgium. Association for Computational Linguistics.
- Xiaodan Lee, Yang Yao, Yaowei Zhang, Yanyan Guan, and Xiaolin Rui. 2015. Emotion classification using massive parallel data mined from social media. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Anand Kumar Madasamy and Soman Kutti Padannayil. 2021. Transfer learning based code-mixed part-of-speech tagging using character level representations for indian languages. *Journal of Ambient Intelligence and Humanized Computing*, pages 1–12.
- Asrita Venkata Mandalam and Yashvardhan Sharma. 2021. [Sentiment analysis of Dravidian code mixed data](#). In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 46–54, Kyiv. Association for Computational Linguistics.
- Shubhanshu Mishra and Sudhanshu Mishra. 2019. 3idiots at hasoc 2019: Fine-tuning transformer neural networks for hate speech identification in indo-european languages. In *FIRE (Working Notes)*, pages 208–213.
- P. Mæhlum, J. Barnes, L. Øvrelid, and E. Velldal. 2019. Annotating evaluative sentences for sentiment analysis: a dataset for norwegian. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, pages 121–130, Turku, Finland. Linköping University Electronic Press.
- Sridevi Padmamala and R Vijayarani. 2017. Sentiment analysis in tamil: A comparative study. In *Proceedings of the International Conference on Computing, Communication and Automation (ICCCA)*.
- Sharmistha Patra, Monojit Choudhury, and Animesh Mukherjee. 2018. Sentiment analysis in code-mixed social media text. In *Proceedings of the 11th International Conference on Language Resources and Evaluation (LREC)*.
- Prasanth, R Aswin Raj, Adhithan P, Premjith B, and Soman Kp. 2022. [CEN-Tamil@DravidianLangTech-ACL2022: Abusive comment detection in Tamil using TF-IDF and random kitchen sink algorithm](#). In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 70–74, Dublin, Ireland. Association for Computational Linguistics.
- Ruba Priyadharshini, Bharathi Raja Chakravarthi, Malliga S, SUBALALITHA CN, Kogilavani S V, Premjith B, Abirami Murugappan, and Prasanna Kumar" Kumaresan. 2023. Overview of shared-task on abusive comment detection in tamil and telugu. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.
- Satish Ranjan, Sandeep Kumar, Monojit Choudhury, and Shivakumar Patel. 2016. A comparative study of sentiment analysis in hindi, english, and code-mixed data. In *Proceedings of the 8th Indian Conference on Human Computer Interaction (HCI)*.
- A. Rogers, A. Romanov, A. Rumshisky, S. Volkova, M. Gronas, and A. Gribov. 2018. Rusentiment: An enriched sentiment analysis dataset for social media in russian. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 755–763, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Sunayana Sitaram, AR Balamurali, and Shreekantha Rakshit. 2015. Sentiment analysis of code-mixed tweets. In *Proceedings of the International Conference on Advances in Computing, Communications and Informatics (ICACCI)*.
- Thamar Solorio, Elizabeth Blair, and Suraj Maharjan. 2014. Overview for the first shared task on language identification in code-switched data. In *Proceedings of the First Workshop on Computational Approaches to Code Switching*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz

- Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#).
- David Vilares, Miguel A Alonso, Carlos Gómez-Rodríguez, Helena Gómez-Adorno, and Thamar Solorio. 2015. Sentiment analysis on monolingual, multilingual and code-switching twitter corpora. In *Proceedings of the 16th Annual Conference of the International Speech Communication Association (INTERSPEECH)*.
- David Vilares, Miguel A Alonso, Carlos Gómez-Rodríguez, Helena Gómez-Adorno, and Thamar Solorio. 2016. Semeval-2016 task 4: Sentiment analysis in twitter. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval)*.
- J. Wiebe, T. Wilson, and C. Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language Resources and Evaluation*, 39(2):165–210.
- G. I. Winata, Z. Lin, and P. Fung. 2019a. Learning multilingual meta-embeddings for code-switching named entity recognition. In *Proceedings of the 4th Workshop on Representation Learning for NLP (RepL4NLP-2019)*, pages 181–186, Florence, Italy. Association for Computational Linguistics.
- Genta Indra Winata, Yik-Cheung Lim, and Erik Cambria. 2019b. Code-switched sentiment analysis with pretrained contextualized embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Genta Indra Winata, Yik-Cheung Lim, and Erik Cambria. 2019c. Hierarchical meta-embeddings for code-mixed sentiment analysis. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Stefan Ziehe, Franziska Pannach, and Aravind Krishnan. 2021. [GCDH@LT-EDI-EACL2021: XLM-RoBERTa for hope speech detection in English, Malayalam, and Tamil](#). In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 132–135, Kyiv. Association for Computational Linguistics.