

GJG@TamilNLP-ACL2022: Using Transformers for Abusive Comment Classification in Tamil

Gaurang Prasad*

wikiHow Inc.
Palo Alto, CA, USA
gaurang@wikihow.com

Janvi Prasad*

Vellore Institute of Technology
Vellore, India
janvi.prasad@gmail.com

Gunavathi Chellamuthu

Vellore Institute of Technology
Vellore, India
gunavathi.cm@vit.ac.in

Abstract

This paper presents transformer-based models for the "Abusive Comment Detection" shared task at the Second Workshop on Speech and Language Technologies for Dravidian Languages at ACL 2022. Our team participated in both the multi-class classification sub-tasks as a part of this shared task. The dataset for sub-task A was in Tamil text; while B was code-mixed Tamil-English text. Both the datasets contained 8 classes of abusive comments. We trained an XLM-RoBERTa and DeBERTa base model on the training splits for each sub-task. For sub-task A, the XLM-RoBERTa model achieved an accuracy of 0.66 and the DeBERTa model achieved an accuracy of 0.62. For sub-task B, both the models achieved a classification accuracy of 0.72; however, the DeBERTa model performed better in other classification metrics. Our team ranked 2nd in the code-mixed classification sub-task and 8th in Tamil-text sub-task.

1 Introduction

The advent of social media and social networks have completely revolutionized the way people communicate with one another (Priyadharshini et al., 2021; Kumaresan et al., 2021). There are many positive aspects of social media - improved connectivity, real-time conversation across multiple locations, a new type of social construct, etc. However, this surge of internet-based communication has also brought about an increase in the volume of negative comments (Sampath et al., 2022; Ravikiran et al., 2022; Chakravarthi et al., 2022; Bharathi et al., 2022; Priyadharshini et al., 2022). Being able to detect and classify such negative and abusive comments is a fundamental and challenging problem to solve.

It is fundamental because better hate & abusive comment detection leads to the improvement of spam detection systems, improves web-inclusivity, and ultimately makes the internet a better place for everybody (Ghanghor et al., 2021a,b; Yasaswini et al., 2021).

Abusive Comment detection and classification falls under the broader spectrum of text classification tasks in Natural Language Processing (NLP). Improvements in abusive comment classifiers can directly enhance content filtering systems, digital well-being software, spam detection, etc (Chakravarthi et al., 2021b, 2020). It also leads to a better physical and mental experience for the end-user, as these classifiers can be used to reduce the various types of abusive comments in the digital world.

Just as with any other NLP task, building good abusive comment classifiers requires annotated datasets. There are plenty of such datasets available for high-resource languages like English, which has led to a lot of published research in this space. However, for low-resource languages like Tamil, there are very few publicly available datasets for downstream NLP tasks (Anita and Subalalitha, 2019b,a; Subalalitha and Poovammal, 2018). Tamil is a Dravidian classical language used by the Tamil people of South Asia. Tamil is an official language of Tamil Nadu, Sri Lanka, Singapore, and the Union Territory of Puducherry in India (Subalalitha, 2019; Srinivasan and Subalalitha, 2019; Narasimhan et al., 2018). Significant minority speak Tamil in the four other South Indian states of Kerala, Karnataka, Andhra Pradesh, and Telangana, as well as the Union Territory of the Andaman and Nicobar Islands. It is also spoken by the Tamil diaspora, which may be found in Malaysia, Myanmar, South Africa, the United Kingdom, the United States, Canada, Australia,

*These authors contributed equally to this work

and Mauritius (Sakuntharaj and Mahesan, 2021, 2017, 2016; Thavareesan and Mahesan, 2019, 2020a,b, 2021). Tamil is also the native language of Sri Lankan Moors. Tamil, one of the 22 scheduled languages in the Indian Constitution, was the first to be designated as a classical language of India (B and A, 2021b,a).

This shared task is an attempt to promote research in abusive comment detection and classification in Tamil. This novel dataset, generated from YouTube comments, consists of 8 types of abusive comments. The publication of the Tamil-text as well as the code-mixed Tamil-English datasets also provides an opportunity to learn about the performance of models on different character sets. In this shared task, we participated in both the sub-tasks: the Tamil-text classification sub-task A, and the code-mixed sub-task B. The goal of this paper is to demonstrate the performance of fine-tuning pre-trained transformer-based models for such a text-classification task in Tamil. We train an XLM-RoBERTa and DeBERTa model for each sub-task on the given train splits, optimize parameters, and evaluate the performance on the test split.

The rest of our paper is organized as follows: we discuss related work in Tamil emotion recognition, describe the datasets, our methodology, and conclude with the results and performance metrics.

We provide a link to our models and evaluations to provide reproducibility, and empower future research in this domain¹.

2 Related Work

As mentioned above, a lot of published work exists in the domain of offensive language detection for high-resource languages. Kwok and Wang (2013) trained a binary classifier to detect racist tweets in English. Xu et al. (2012) presented off-the-shelf NLP approaches to identify bullying in social media, in English. Kumar et al. (2018) present their findings from a shared task for aggression identification in social media. Nobata et al. (2016) demonstrated a Machine Learning approach to detect abusive language in English online content.

The volume of published research in abusive comment detection for low-resource languages is much lower than that of English and other high-resource languages. Wiegand et al. (2018) provided an overview of a shared abusive comment detection task in German. Kannan and Mitrović

(2021), Kamal et al. (2021), and Jha et al. (2020) have presented their work in detecting abusive comments in Hindi. Eshan and Hasan (2017), Emon et al. (2019), and Romim et al. (2021) have demonstrated popular approaches for Bengali.

Chakravarthi (2020), Chakravarthi and Muralidaran (2021), and Hande et al. (2021) have published their findings from other text classification-related shared tasks in Dravidian Languages (hope speech detection).

Mandl et al. (2020) organized a workshop track for hate speech detection in Tamil, Malayalam, Hindi, English and German. We believe this was the first big focus on developing abusive content identification and classification techniques for Dravidian languages. This was followed up by Chakravarthi et al. (2021a), who organized a shared task for offensive language identification in Tamil, Malayalam, and Kannada.

There have also been multiple previous works that use XLM-RoBERTa for text classification tasks. Zhao and Tao (2021) proposed a system using XLM-RoBERTa and DPCNN for detecting offensive text in Dravidian languages. Qu et al. (2021) used TextCNN and XLM-RoBERTa from emotion classification in Spanish. Ou and Li (2020) also demonstrated using XLM-RoBERTa for a hate speech identification classification task.

The number of published works using DeBERTa is fewer than that of XLM-RoBERTa. There have been some studies that use DeBERTa for entity extraction and text-classification tasks. (Martin and Pedersen, 2021; Khan et al., 2022)

3 Data

The organizers of the shared task released the annotated training and development splits for both the sub-tasks. The testing dataset, without labels, was released a few days prior to the run submission deadline. Once the results were announced, the organizers released the labeled test dataset for verification purposes.

For sub-task A, the dataset consisted of Tamil sentences annotated to one of eight English abusive categories: Misandry, Counter Speech (Sp.), Misogyny, Xenophobia, Hope Sp., Homophobia, Transphobia, or None of the Above (N.O.T.A). The dataset for sub-task B was code-mixed Tamil-English with the same eight English abusive comment classes as sub-task A. The average length of a sentence in the corpora was 1 (Priyadharshini

¹<https://tinyurl.com/GJGAbusiveComments>

Label	Count	Label	Count
N.O.T.A	1296	Hope Sp.	86
Misandry	446	Homophob.	35
Counter Sp.	149	Transphob.	6
Misogyny	125	Not Tamil	2
Xenophobia	95		

Table 1: Classification labels for sub-task A and the number of rows under each label in the train split.

Label	Count	Label	Count
N.O.T.A	3720	Hope Sp.	213
Misandry	830	Misogyny	211
Counter Sp.	348	Homophob.	172
Xenophobia	297	Transphob.	157

Table 2: Classification labels for sub-task B and the number of rows under each label in the training split.

et al., 2022). The sentences in both the datasets contained extended words, special characters, emojis, grammatical, and lexical inconsistencies.

3.1 Sub-Task A

The training split for classification sub-task A had a total of 2,240 rows. 2,238 of these rows were classified under one of the eight categories of abusive comments. 2 rows were not in Tamil. Including the "Not Tamil" category, there were a total of 9 category labels in the test split. The development dataset contained 560 rows. The test split had 698 rows. The classification labels along with the total count for each label in the train split are represented in Table 1.

3.2 Sub-Task B

Sub-task B used the code-mixed Tamil-English dataset with the same 8 abusive category labels. The training split did not include any rows of the "Not Tamil" category, which was seen in sub-task A. The dataset splits were larger than that of sub-task A: the training split had a total of 5,948 rows. The development dataset had 1,488 rows and the test split had 1,856 rows. Table 2 represents all the class labels in the train split along with the total number of rows under each label.

4 Methodology

For each classification sub-task, we fine-tune an XLM-RoBERTa and DeBERTa base model on sentences from the training splits to create a classification model. We do not remove any stop

words, special characters, or emojis from the test splits, in order to preserve the context of the comment. Special characters and emojis provide useful context, especially for a text classification task. For sub-task A, we also did not remove the 2 instances of "Not Tamil" from the training set.

XLM-RoBERTa is a multilingual version of RoBERTa, which in itself was an improvement over BERT to achieve state-of-the-art results in multiple NLP tasks. XLM-RoBERTa is pre-trained on 2.5TB of filtered CommonCrawl data containing 100 languages (Conneau et al., 2019). DeBERTa uses disentangled attention and enhanced mask decoder to enhance RoBERTa and outperform it in a majority of NLP tasks (He et al., 2021).

Table 3 represents the parameters used to fine-tune the XLM-RoBERTa and DeBERTa base models for both the sub-tasks.

5 Results

We use accuracy and the weighted averages of precision, recall, and F1-score as performance metrics to evaluate our classification models. We calculate all four evaluation metrics to get a better sense of the classification performance.

For sub-task A, we find that the XLM-RoBERTa outperforms the DeBERTa in all evaluation metrics. Both models used the same training splits and parameters. The multi-lingual nature of XLM-RoBERTa is evident from its better performance.

The overall classification performance for sub-task 2 was higher than that for sub-task A. We believe this is because of the English character set used by the code-mixed dataset. Both the transformer-models are pre-trained on large English datasets. DeBERTa outperforms XLM-RoBERTa in all metrics, which justifies DeBERTa's improvement over XLM-RoBERTa for English-character-NLP tasks (Table 4).

6 Conclusion and Future Work

We present XLM-RoBERTa and DeBERTa models for two multi-class text classification tasks in Tamil. The objective of the shared task was to identify abusive content in Tamil text. There were two sub-tasks: sub-task A in Tamil text, and sub-task B which used Tamil-English code-mixed text. The classes of abusive comments were the same for both the sub-tasks. The Tamil dataset, for sub-task A, consisted of 2,240 rows in the training split. The

	Parameter Value			
	Sub-Task A		Sub-Task B	
	XLM-RoBERTa	DeBERTa	XLM-RoBERTa	DeBERTa
Batch Size	9	9	10	10
Max. Sequence Length	256	256	256	256
Number of Epochs	10	10	10	10
Learning Rate	1e-5	1e-5	1e-5	1e-5
Weight Decay	0	0	0	0
Use Class Weights	False	False	False	False

Table 3: Fine-tuning parameters of XLM-RoBERTa and DeBERTa models for both sub-tasks

Task	Model	Accuracy	F1-score	Precision	Recall
Sub-Task A	XLM-RoBERTa	0.66	0.65	0.65	0.66
	DeBERTa	0.62	0.57	0.62	0.56
Sub-Task B	XLM-RoBERTa	0.72	0.70	0.70	0.72
	DeBERTa	0.72	0.72	0.72	0.72

Table 4: Classification metrics for both sub-tasks.

code-mixed training split, for sub-task B, was much bigger with 5,948 rows.

We propose the fine-tuning of pre-trained XLM-RoBERTa and DeBERTa, which are transformer-based models, for classifying Tamil text into abusive comment classes. We trained all the models using the respective training splits as-is. For sub-task A, the XLM-RoBERTa achieved a classification accuracy of 66% with a weighted F-1 of 0.65, precision of 0.65, and a recall value of 0.66. The DeBERTa model achieved an accuracy of 62% with weighted F-1 of 0.57, precision of 0.62, and 0.56 recall. For sub-task B, the XLM-RoBERTa achieved a classification accuracy of 72% with a weighted F-1 of 0.70, precision of 0.70, and a recall value of 0.72. The DeBERTa model achieved an accuracy of 72% with weighted F-1 of 0.72, precision of 0.72, and 0.72 recall.

We show that the XLM-RoBERTa model outperforms DeBERTa for sub-task A, which used Tamil text. For the code-mixed Tamil-English text (sub-task B), the DeBERTa model outperforms the XLM-RoBERTa. This validates the strength of the DeBERTa model on English character tasks, and the superiority of the XLM-RoBERTa for non-English languages. By using the training split as-is, we retain the information provided by special characters like emojis and extended punctuation symbols. The XLM-RoBERTa model had the eight best classification performance in the shared task for sub-task A, and the DeBERTa model ranked

second best. We have open-sourced the code used in this study in a public GitHub repository.

For future-work, removing the "Not Tamil" rows from the training split for sub-task A would eliminate an extremely under-sampled class, which may lead to performance improvements. Extracting emojis separately from the text and incorporating emoji-information in identifying abusive comments is also an area for potential study. Using pre-trained models is a good starting point in this domain, however, we feel that custom model architectures and systems have to be studied for such classification tasks in Tamil.

References

- R Anita and CN Subalalitha. 2019a. An approach to cluster Tamil literatures using discourse connectives. In *2019 IEEE 1st International Conference on Energy, Systems and Information Processing (ICESIP)*, pages 1–4. IEEE.
- R Anita and CN Subalalitha. 2019b. Building discourse parser for Thirukkural. In *Proceedings of the 16th International Conference on Natural Language Processing*, pages 18–25.
- Bharathi B and Agnusimmaculate Silvia A. 2021a. [SSNCSE_NLP@DravidianLangTech-EACL2021: Meme classification for Tamil using machine learning approach](#). In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 336–339, Kyiv. Association for Computational Linguistics.
- Bharathi B and Agnusimmaculate Silvia A. 2021b.

- SSNCSE_NLP@DravidianLangTech-EACL2021: Offensive language identification on multilingual code mixing text. In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 313–318, Kyiv. Association for Computational Linguistics.
- B Bharathi, Bharathi Raja Chakravarthi, Subalalitha Chinnaudayar Navaneethakrishnan, N Sripriya, Arunaggiri Pandian, and Swetha Valli. 2022. Findings of the shared task on Speech Recognition for Vulnerable Individuals in Tamil. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*. Association for Computational Linguistics.
- Bharathi Raja Chakravarthi. 2020. [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*, pages 41–53, Barcelona, Spain (Online). Association for Computational Linguistics.
- Bharathi Raja Chakravarthi and Vigneshwaran Muralidaran. 2021. [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*, pages 61–72, Kyiv. Association for Computational Linguistics.
- Bharathi Raja Chakravarthi, Vigneshwaran Muralidaran, Ruba Priyadharshini, and John Phillip 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.
- Bharathi Raja Chakravarthi, Ruba Priyadharshini, Thenmozhi Durairaj, John Phillip McCrae, Paul Buitaleer, Prasanna Kumar Kumaresan, and Rahul Ponnusamy. 2022. Findings of the shared task on Homophobia Transphobia Detection in Social Media Comments. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*. Association for Computational Linguistics.
- Bharathi Raja Chakravarthi, Ruba Priyadharshini, Navya Jose, Thomas Mandl, Prasanna Kumar Kumaresan, Rahul Ponnusamy, RL Hariharan, John Phillip McCrae, Elizabeth Sherly, et al. 2021a. 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*, pages 133–145.
- Bharathi Raja Chakravarthi, Ruba Priyadharshini, Rahul Ponnusamy, Prasanna Kumar Kumaresan, Kayalvizhi Sampath, Durairaj Thenmozhi, Sathiyaraj Thangasamy, Rajendran Nallathambi, and John Phillip McCrae. 2021b. Dataset for identification of homophobia and transphobia in multilingual YouTube comments. *arXiv preprint arXiv:2109.00227*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Unsupervised cross-lingual representation learning at scale](#). *CoRR*, abs/1911.02116.
- Estiak Ahmed Emon, Shihab Rahman, Joti Banarjee, Amit Kumar Das, and Tanni Mitra. 2019. A deep learning approach to detect abusive bengali text. In *2019 7th International Conference on Smart Computing & Communications (ICSCC)*, pages 1–5. IEEE.
- Shahnoor C Eshan and Mohammad S Hasan. 2017. An application of machine learning to detect abusive bengali text. In *2017 20th International Conference of Computer and Information Technology (ICCIT)*, pages 1–6. IEEE.
- Nikhil Ghanghor, Parameswari Krishnamurthy, Sajeetha Thavareesan, Ruba Priyadharshini, and Bharathi Raja Chakravarthi. 2021a. [IITK@DravidianLangTech-EACL2021: Offensive language identification and meme classification in Tamil, Malayalam and Kannada](#). In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 222–229, Kyiv. Association for Computational Linguistics.
- Nikhil Ghanghor, Rahul Ponnusamy, Prasanna Kumar Kumaresan, Ruba Priyadharshini, Sajeetha Thavareesan, and Bharathi Raja Chakravarthi. 2021b. [IITK@LT-EDI-EACL2021: Hope speech detection for equality, diversity, and inclusion in Tamil, Malayalam and English](#). In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 197–203, Kyiv. Association for Computational Linguistics.
- Adeep Hande, Ruba Priyadharshini, Anbukkarasi Sampath, Kingston Pal Thamburaj, Prabakaran Chandran, and Bharathi Raja Chakravarthi. 2021. [Hope speech detection in under-resourced kannada language](#).
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [Deberta: Decoding-enhanced bert with disentangled attention](#). In *International Conference on Learning Representations*.
- Vikas Kumar Jha, Pa Hrudya, PN Vinu, Vishnu Vijayan, and Pa Prabakaran. 2020. Dhoot-repository and classification of offensive tweets in the hindi language. *Procedia Computer Science*, 171:2324–2333.

- Ojasv Kamal, Adarsh Kumar, and Tejas Vaidhya. 2021. Hostility detection in hindi leveraging pre-trained language models. In *International Workshop on Combating On line Hostile Posts in Regional Languages during Emergency Situation*, pages 213–223. Springer.
- Sudharsana Kannan and Jelena Mitrović. 2021. Hatespeech and offensive content detection in hindi language using c-bigru. In *Forum for Information Retrieval Evaluation (Working Notes)(FIRE)*, CEUR-WS.org.
- Pervaiz Iqbal Khan, Imran Razzak, Andreas Dengel, and Sheraz Ahmed. 2022. Performance comparison of transformer-based models on twitter health mention classification. *IEEE Transactions on Computational Social Systems*.
- 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)*, pages 1–11.
- Prasanna Kumar Kumaresan, Ratnasingam Sakuntharaj, Sajeetha Thavareesan, Subalalitha Navaneethakrishnan, Anand Kumar Madasamy, Bharathi Raja Chakravarthi, and John P McCrae. 2021. Findings of shared task on offensive language identification in Tamil and Malayalam. In *Forum for Information Retrieval Evaluation*, pages 16–18.
- Irene Kwok and Yuzhou Wang. 2013. Locate the hate: Detecting tweets against blacks. In *Twenty-seventh AAAI conference on artificial intelligence*.
- 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*, pages 29–32.
- Anna Martin and Ted Pedersen. 2021. Duluth at semeval-2021 task 11: Applying deberta to contributing sentence selection and dependency parsing for entity extraction. In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 490–501.
- Anitha Narasimhan, Aarthi Anandan, Madhan Karky, and CN Subalalitha. 2018. Porul: Option generation and selection and scoring algorithms for a tamil flash card game. *International Journal of Cognitive and Language Sciences*, 12(2):225–228.
- Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. 2016. Abusive language detection in online user content. In *Proceedings of the 25th international conference on world wide web*, pages 145–153.
- Xiaozhi Ou and Hongling Li. 2020. Ynu_oxz@haspeede 2 and ami: Xlm-roberta with ordered neurons lstm for classification task at evalita 2020. *EVALITA Evaluation of NLP and Speech Tools for Italian*, 2765:102–109.
- Ruba Priyadharshini, Bharathi Raja Chakravarthi, Subalalitha Chinnaudayar Navaneethakrishnan, Thenmozhi Durairaj, Malliga Subramanian, Kogilavani Shanmugavadivel, Siddhanth U Hegde, and Prasanna Kumar Kumaresan. 2022. Findings of the shared task on Abusive Comment Detection in Tamil. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*. Association for Computational Linguistics.
- Ruba Priyadharshini, Bharathi Raja Chakravarthi, Sajeetha Thavareesan, Dhivya Chinnappa, Durairaj Thenmozhi, and Rahul Ponnusamy. 2021. Overview of the DravidianCodeMix 2021 shared task on sentiment detection in Tamil, Malayalam, and Kannada. In *Forum for Information Retrieval Evaluation*, pages 4–6.
- S Qu, Y Yang, and Q Que. 2021. Emotion classification for spanish with xlm-roberta and textcnn. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2021)*. CEUR Workshop Proceedings, CEUR-WS, Málaga, Spain.
- Manikandan Ravikiran, Bharathi Raja Chakravarthi, Anand Kumar Madasamy, Sangeetha Sivanesan, Ratnavel Rajalakshmi, Sajeetha Thavareesan, Rahul Ponnusamy, and Shankar Mahadevan. 2022. Findings of the shared task on Offensive Span Identification in code-mixed Tamil-English comments. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*. Association for Computational Linguistics.
- Nauros Romim, Mosahed Ahmed, Hriteshwar Talukder, Saiful Islam, et al. 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*, pages 457–468. Springer.
- Ratnasingam Sakuntharaj and Sinnathamby Mahesan. 2016. A novel hybrid approach to detect and correct spelling in Tamil text. In *2016 IEEE International Conference on Information and Automation for Sustainability (ICIAfS)*, pages 1–6.
- Ratnasingam Sakuntharaj and Sinnathamby Mahesan. 2017. Use of a novel hash-table for speeding-up suggestions for misspelt Tamil words. In *2017 IEEE International Conference on Industrial and Information Systems (ICIIS)*, pages 1–5.
- Ratnasingam Sakuntharaj and Sinnathamby Mahesan. 2021. Missing word detection and correction based on context of Tamil sentences using n-grams. In *2021 10th International Conference on Information and Automation for Sustainability (ICIAfS)*, pages 42–47.

- Anbukkarasi Sampath, Thenmozhi Durairaj, Bharathi Raja Chakravarthi, Ruba Priyadharshini, Subalalitha Chinnaudayar Navaneethkrishnan, Kogilavani Shanmugavadivel, Sajeetha Thavareesan, Sathiyaraj Thangasamy, Parameswari Krishnamurthy, Adeep Hande, Sean Benhur, and Santhiya Ponnusamy, Kishor Kumar Pandiyan. 2022. Findings of the shared task on Emotion Analysis in Tamil. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*. Association for Computational Linguistics.
- R Srinivasan and CN Subalalitha. 2019. Automated named entity recognition from tamil documents. In *2019 IEEE 1st International Conference on Energy, Systems and Information Processing (ICESIP)*, pages 1–5. IEEE.
- C. N. Subalalitha. 2019. [Information extraction framework for Kurunthogai](#). *Sādhanā*, 44(7):156.
- CN Subalalitha and E Poovammal. 2018. Automatic bilingual dictionary construction for Tirukurul. *Applied Artificial Intelligence*, 32(6):558–567.
- Sajeetha Thavareesan and Sinnathamby Mahesan. 2019. [Sentiment analysis in Tamil texts: A study on machine learning techniques and feature representation](#). In *2019 14th Conference on Industrial and Information Systems (ICIIS)*, pages 320–325.
- Sajeetha Thavareesan and Sinnathamby Mahesan. 2020a. [Sentiment lexicon expansion using Word2vec and fastText for sentiment prediction in Tamil texts](#). In *2020 Moratuwa Engineering Research Conference (MERCon)*, pages 272–276.
- Sajeetha Thavareesan and Sinnathamby Mahesan. 2020b. [Word embedding-based part of speech tagging in Tamil texts](#). In *2020 IEEE 15th International Conference on Industrial and Information Systems (ICIIS)*, pages 478–482.
- Sajeetha Thavareesan and Sinnathamby Mahesan. 2021. [Sentiment analysis in Tamil texts using k-means and k-nearest neighbour](#). In *2021 10th International Conference on Information and Automation for Sustainability (ICIAfS)*, pages 48–53.
- Michael Wiegand, Melanie Siegel, and Josef Ruppenhofer. 2018. Overview of the germeval 2018 shared task on the identification of offensive language.
- Jun-Ming Xu, Kwang-Sung Jun, Xiaojin Zhu, and Amy Bellmore. 2012. Learning from bullying traces in social media. In *Proceedings of the 2012 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 656–666.
- Konthala Yasaswini, Karthik Puranik, Adeep Hande, Ruba Priyadharshini, Sajeetha Thavareesan, and Bharathi Raja Chakravarthi. 2021. [IIT@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*, pages 187–194, Kyiv. Association for Computational Linguistics.
- Yingjia Zhao and Xin Tao. 2021. [Zyj123@dravidianlangtech-eacl2021: Offensive language identification based on xlm-roberta with dpcnn](#). In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 216–221.