# Overview of Shared-task on Abusive Comment Detection in Tamil and Telugu

**Ruba Priyadharshini[1], Bharathi Raja Chakravarthi[2],**
**Malliga Subramanian[3], Subalalitha Chinnaudayar Navaneethakrishnan[4],**
**Kogilavani Shanmugavadivel[3], Premjith B[5], Abirami Murugappan[6],**
**Prasanna Kumar Kumaresan[2], Karnati Sai Prashanth[5],**
**Mangamuru Sai Rishith Reddy[5], Janakiram Chandu[5]**

[1]Gandhigram Rural Institute-Deemed to be University, India
[2]Insight SFI Research Centre for Data Analytics, School of Computer Science,
University of Galway, Ireland
[3]Kongu Engineering College, Tamil Nadu, India
[4]SRM Institute of Science and Technology, Kattankulathur, Chennai, India
[5]Amrita School of Artificial Intelligence, Amrita Vishwa Vidyapeetham, India
[6]Department of Information Science and Technology, Anna University

## Abstract

This paper discusses the submissions to the shared task on abusive comment detection in Tamil and Telugu codemixed social media text conducted as part of the third Workshop on Speech and Language Technologies for Dravidian Languages at RANLP 2023. The task encourages researchers to develop models to detect the contents containing abusive information in Tamil and Telugu codemixed social media text. The task has three subtasks - abusive comment detection in Tamil, Tamil-English and Telugu-English. The dataset for all the tasks was developed by collecting comments from YouTube. The submitted models were evaluated using macro F1-score, and prepared the rank list accordingly.

## 1 Introduction

Abusive comment detection from social media has become an essential and challenging task in the current age of technology (Chakravarthi et al., 2023; Priyadharshini et al., 2022b; Prasanth et al., 2022). The proliferation of online platforms helped people to spread information, including harmful and violent comments and posts. Therefore, addressing and mitigating harmful content to keep online platforms clean automatically has become very important (Chakravarthi et al., 2022a,b, 2023; Chakravarthi, 2023). This task is challenging due to the complexities of the languages. However, advanced machine learning algorithms and techniques were proposed to automatically identify and flag abusive comments, ranging from hate speech and cyberbullying to threats and harassment, recently. These systems analyze the content of the posts and the context to determine the presence of abusive language and malicious intent. The complexity of detecting the abusive contents from a code-mixed Dravidian language is even high due to the code-mixed nature of the text and the intricacies of the language, such as morphological richness and agglutinative property (Premjith et al., 2018). In addition, large datasets of labelled abusive content are required to train and fine-tune the Artificial Intelligence (AI)-based models, enabling them to recognize patterns and distinguish between harmful and benign texts.

A considerable amount of research has been conducted to detect abusive and similar harmful content from social media posts and comments (Bharathi and Agnusimmaculate Silvia, 2021; Bharathi and Varsha, 2022b; Swaminathan et al., 2022; Subramanian et al., 2022). In addition, several shared tasks were organized to promote the research for automatically detecting social media comments containing abusive content. This shared task focuses on detecting abusive comments in two Dravidian languages - Tamil and Telugu. Tamil is predominantly spoken in Tamil Nadu, a state in India and nearby countries, whereas Telugu is the official language of two states in India - Andhra Pradesh and Telangana (Vasantharajan et al., 2022; Anita and Subalalitha, 2019; Thavareesan and Mahesan, 2019, 2020a,b; Subalalitha, 2019; Sakuntharaj and Mahesan, 2016, 2017, 2021). This paper summarizes the findings of the research works submitted to the shared task on abusive language detection in Tamil and Telugu. Besides, this paper details the dataset developed and used for conducting the experiments.

## 2 Literature Review

In an attempt by Chakravarthy et. al. (Chakravarthi et al., 2023), a set of four datasets comprising abusive comments in Tamil and code-mixed Tamil-English extracted from YouTube is presented. Each dataset has undergone comment-level annotation, wherein polarities are assigned to the comments. To establish baselines for these datasets, the authors conducted experiments using various machine learning classifiers and presented the results in terms of F1-score, precision, and recall. Prasanth et al. (Prasanth et al., 2022) conducted a study focusing on the detection of abusive comments within given text. The authors employed TF-IDF with char-wb analyzers and utilized the Random Kitchen Sink (RKS) algorithm to generate feature vectors. For classification purposes, they employed the Support Vector Machine (SVM) classifier with a polynomial kernel. The proposed method was applied to both Tamil and Tamil-English datasets, resulting in f1-scores of 0.32 and 0.25, respectively. Priyadharsini et. al. (Priyadharshini et al., 2022a) provides a comprehensive review of a shared task focused on identifying abusive comments encompassing various categories such as Homophobia, Misandry, Counter-speech, Misogyny, Xenophobia, Transphobic, and hope speech. The participants were provided with a dataset extracted from social media, which was labeled with the aforementioned categories in both Tamil and Tamil-English code-mixed languages. The participants employed diverse machine learning and deep learning algorithms for their approaches. The paper presents an overview of this task, including detailed information about the dataset and the results achieved by the participants. The objective of the study by Bharathi and Varsha (Bharathi and Varsha, 2022a) is to automate the identification and categorization of abusive comments into specific categories like Misogynism, Misandry, Homophobia, and Cyberbullying. The datasets utilized in this research were provided by the DravidianLangTech@ACL2022 organizers and consisted of code-mixed Tamil text. The authors trained these datasets using pre-trained transformer models such as BERT, m-BERT, and XLNET. Remarkably, they achieved a weighted average F1 score of 0.96 for Tamil-English code-mixed text and 0.59 for Tamil text. Gupta et al. (Gupta et al., 2022) conducted a study where they introduced a model called AbuseXLMR, designed specifically for detecting abusive content. This model was pre-trained on a vast amount of social media comments in over 15 Indic languages. Notably, AbuseXLMR exhibited superior performance compared to XLM-R and MuRIL when evaluated on multiple Indic datasets. In addition to providing annotations, this study also released mappings between comment, post, and user IDs, enabling the modeling of relationships among them. Furthermore, competitive baselines for monolingual, cross-lingual, and few-shot scenarios were shared, intending to establish the collected dataset as a benchmark for future research. The primary goal of the study by Marreddy et. al. (Marreddy et al., 2022) is to address the challenges posed by limited resources in the Telugu language. The authors make several valuable contributions to enrich resources for Telugu. They have curated a large annotated dataset containing 35,142 sentences for various NLP tasks, including sentiment analysis, emotion identification, hate-speech detection, and sarcasm detection. To enhance model efficiency, the authors have developed separate lexicons for sentiment, emotion, and hate speech and utilized pre-trained word and sentence embeddings. Furthermore, the authors have created different pre-trained language models specifically for Telugu, such as ELMo-Te, BERT-Te, RoBERTa-Te, ALBERT-Te, and DistilBERT-Te, using a sizable Telugu corpus comprising 8,015,588 sentences. Notably, the authors demonstrate that these developed models significantly enhance the performance of the four NLP tasks and provide benchmark results for Telugu.

## 3 Task Description

We used the CodaLab platform to conduct the task [1]. The task includes three subtasks - abusive language detection in

- **Tamil**: Abusive language detection from Tamil codemixed social media text

- **Tamil-English**: Abusive language detection from Tamil-English codemixed social media text

- **Telugu-English**: Abusive language detection from Telugu-English codemixed social media text

### 3.1 Tamil

The dataset was compiled utilizing the YouTube comment scraper, capturing comprehensive com-

---

[1]https://codalab.lisn.upsaclay.fr/competitions/11096

ments in the Tamil script. The comments were sourced from videos addressing subjects related to homophobia, transphobia, misogyny, xenophobia, and misandry. However, procuring Tamil comments from YouTube videos posed challenges due to the extensive array of videos available. An extensive effort was made to exclusively retain comments in the Tamil language, resulting in the exclusion of non-Tamil comments. The comment annotations encompassed seven distinct classes, which are itemized in Table 1. The table provides comment counts for each class within every set.

Table 1: Distribution of training, test, and dev datasets used for the shared task on abusive language detection in Tamil

| Categories | Train | Test | Dev |
|---|---|---|---|
| None-of-the-above | 1295 | 416 | 346 |
| Hope-Speech | 86 | 26 | 11 |
| Homophobia | 35 | 8 | 8 |
| Misandry | 446 | 127 | 104 |
| Counter-speech | 149 | 47 | 36 |
| Transphobic | 9 | 2 | 2 |
| Xenophobia | 95 | 25 | 29 |
| Misogyny | 125 | 48 | 24 |
| **Total** | **2240** | **699** | **560** |

### 3.2 Tamil-English

The dataset was acquired from YouTube using the YouTube comment scraper. These comments are specifically in Tamil-English codemixed social media text, where Tamil characters are transliterated into the Latin script. Comments that incorporate both Tamil and English words, written in their respective scripts, were included in the dataset. We adhered to YouTube's guidelines to categorize the comments into 7 labels: Homophobia, Transphobia, Hope-speech, Misandry, Xenophobia, Misogyny, Counter-speech, and None-of-the-above. Table 2 below presents the quantities of comments in each dataset as well as the distribution of comments across each label.

### 3.3 Telugu-English

This task was hosted in CodaLab. This task encourages researchers to build machine learning or deep learning models for detecting hate comments from Telugu-English codemixed social media text. The dataset was prepared by collecting hate comments from YouTube. The initial challenge was identifying the videos where we could find the hate com-

Table 2: Distribution of training, test, and dev datasets used for the shared task on abusive language detection in Tamil-English

| Categories | Train | Test | Dev |
|---|---|---|---|
| None-of-the-above | 3720 | 1141 | 919 |
| Hope-Speech | 213 | 70 | 53 |
| Homophobia | 172 | 56 | 43 |
| Misandry | 830 | 292 | 218 |
| Counter-speech | 348 | 88 | 95 |
| Transphobic | 157 | 58 | 40 |
| Xenophobia | 297 | 95 | 70 |
| Misogyny | 211 | 57 | 50 |
| **Total** | **5948** | **1857** | **1488** |

ments. The comments in which Telugu characters are written using Latin scripts and comments containing both Telugu and English words written in respective scripts were considered for preparing the dataset. We followed the regulations of YouTube to annotate the comments into hate and non-hate. The annotators were Telugu native speakers with English proficiency and good academic qualifications. Finally, the dataset consisted of 4500 annotated comments, of which 4000 were used as the training data, and 500 were considered the test data. The training data was released to the participants initially to build the model. The participants were free to choose the validation data. During the testing phase of the competition, we released the test data without labels, and the participants were asked to predict the labels. We published the test data with labels along with the rank list.

The training dataset consisted of 1939 hate and 2061 non-hate comments, whereas the test data had 250 hate and non-hate comments each. The distribution of the data points in each class indicates no considerable class imbalance problem. The train-test split of the dataset and the number of data points in each class is given in Table 3.

Table 3: Distribution of training and test datasets used for the shared task on abusive language detection in Telugu-English

| Category | Train | Test |
|---|---|---|
| Hate | 1939 | 250 |
| Non-hate | 2061 | 250 |
| Total | 4000 | 250 |

We received 52 registrations for the competition. However, only eight teams submitted the predic-

tions for the test data. We accepted a maximum of three runs from each team, and the run with the highest performance score was considered for preparing the rank list, which is shown in Table 5. Macro F1-score was used to evaluate the performance of the submitted results and prepare the rank list.

## 4 System Descriptions

This section summarizes the systems submitted to the Abusive Comment Detection in Tamil and Telugu tasks.

### 4.1 Team: MUCS

The team MUCS (Hegde et al., 2023) submitted three different models to the competition. In all three approaches, the authors used a resampling approach but used different feature extraction algorithms. The first model was developed by using TF-IDF as the feature extraction algorithm. The second method used the Telugu-bert model to generate the input text's feature representation. In contrast, the authors used the multilingual BERT model in the third model. The third approach achieved the highest macro F1-score of 0.7459, and the team secured first place.

### 4.2 Team: DeepBlueAI

The team DeepBlueAI (Luo and Wang, 2023) used XLM-RoBERTa to develop their base model for classifying Telugu comments into hate and non-hate categories. They mixed multiple language datasets at different proportions to build the model. In addition, the authors performed cross-validation to develop a generalized model. This team secured the second position in the shared task, and their submission achieved a macro F1-score of 0.7318.

### 4.3 Team: Habesha

The team followed an LSTM-based approach for modelling the data (Yigezu et al., 2023). They did not use any other algorithms to generate the word embedding. The model consisted of a dropout layer introduced to avoid the overfitting problem, which generally happens when the number of data points is less. In addition, the model was set up to use early stopping based on validation loss, which stops training if the validation loss does not improve for a certain number of epochs. The team was placed third in the competition and scored a macro F1-score of 0.6519.

In the second model, the authors used character-based RNN for training the model.

### 4.4 Team: AK-NLP

The team submitted two models. In the first model, the authors used Term Frequency-Inverse Document Frequency (TF-IDF) to represent the comments as vectors and further used an LSTM model to learn the text's sequential properties. The second approach used a Word2Vec-based hierarchical attention network for building the model. In this work, the authors used 20,000 external codemixed Telugu data for training the Word2Vec model. The TF-IDF+LSTM model achieved the highest performance among the two submissions, with a macro F1-score of 0.6430. This team achieved the fourth rank in the competition.

### 4.5 Team: AbhiPaw

The team AbhiPaw (Bala and Krishnamurthy, 2023) achieved the fifth rank in the compeition, and they scored a macro F1-score of 0.6319. The team implemented a Logistic Regression-based classifier for categorizing the Telugu-Englosh comments into hate and non-hate category.

### 4.6 Team: SuperNova

Team SuperNova (Reddy et al., 2023) used TF-IDF for feature extraction and Support Vector Machine (SVM) for classification. They did not consider any external dataset for feature extraction. The team achieved the sixth position with a macro F1-score of 0.6189.

### 4.7 Team: Athena

The team Athena (Sivanaiah et al., 2023a) implemented their model using the Logistic Regression classifier. They used the TF-IDF vectorization algorithm to vectorize the text data in the dataset before passing it to the model. They did not utilize any external data for generating the feature. The team was placed in the seventh position with a macro F1-score of 0.6137.

### 4.8 Team: CSSCUTN

The team used Bag of Words and TF-IDF feature representation algorithms for converting the input text into a feature representation (Pannerselvam et al., 2023). The authors used machine learning algorithms such as Support Vector Machine (SVM), Logistic Regression, and Random Forest

for building the model. They did not use any external datasets for training the model. The team obtained the eighth

## 5 Result Analysis and Discussion

This section discusses the submission by different teams in the shared task.

### 5.1 Tamil

In the Tamil task, numerous participants took part, contributing a total of 9 submissions. The leading team, MUCS, achieved a macro F1 score of 0.46. This team employed the mBERT pre-trained transformer model using the resample method, along with DistilBERT using the same resample method, both of which delivered the top performance. The second-highest performing team, Harmony, adopted a strategy involving transliteration to Tamil and amalgamation with Tamil data. They balanced class distribution by oversampling minority classes until all classes had an equal count. Additionally, they applied the IndicNLP morphological analyzer for stemming. Subsequently, the data was fed into transformer models: MuRIL and XLM-RoBERTa with fine-tuning. They also employed fast text embedding, which underwent two parallel recurrent layers—two Bi-LSTM and two Bi-GRU and this team got 0.41 macro F1 scores.

### 5.2 Tamil-English

In the Tamil-English codemixed task, 12 participants submitted their evaluation predictions. A rank list was compiled based on the macro F1 scores. Notably, DeepBlueAI secured the top rank with an F1 score of 0.55, while the team Super-Nova achieved the lowest rank with an F1 score of 0.25. The team that claimed the first position employed Fine-tuning with XLM-RoBERTa as the foundational model. They also explored mixing multiple language datasets at various ratios and utilized cross-validation techniques. Conversely, the team that ended up with the last rank adopted a TF-IDF approach in conjunction with basic machine learning models. Interestingly, the teams discovered that the most effective results were achieved when combining Tf-IDF feature extraction with various machine learning and transformer models. Additionally, these teams found success by incorporating resampling techniques into their transformer model implementations.

### 5.3 Telugu-English

The submissions by different teams include various feature extraction approaches and classification models. Most models were based on TF-IDF feature extraction followed by a machine learning classifier. However, the two teams used BERT-based approaches for developing their models. Another two teams used LSTM and RNN architectures for modelling the Telugu-English codemixed data. It is observed from the macro F1-scores of all the teams that the model based on BERT and its variants achieved the top ranks, followed by LSTM and RNN-based models. The bottom-placed teams used TF-IDF feature extraction algorithms. Therefore, it is evident that the BERT-based embedding algorithms learn better features for classification than conventional approaches such as TF-IDF and Bag of Words.

## 6 Conclusion

This paper discussed the findings of the shared task conducted as part of the third Workshop on Speech and Language Technologies for Dravidian Languages at RANLP 2023 on abusive comment detection in Tamil, Tamil-English and Telugu-English data. The datasets used for the competition were collected from YouTube comments and annotated with experts' help in compliance with YouTube's regulations. There were nine, eleven and eight submissions in Tamil, Tamil-English and Telugu-English tasks, respectively. Most teams used multilingual BERT-based pre-trained models to transform the input text into the feature vector. The other submissions consisted of models using TF-IDF features and machine learning classifiers. We used macro F1-score for computing the classification performance and prepared the rank list accordingly.

## References

R Anita and CN Subalalitha. 2019. 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.

Table 4: Rank list for the Tamil subtask

| Team Name | macro F1 | Rank |
|---|---|---|
| MUCS (Hegde et al., 2023) | 0.46 | 1 |
| Harmony (Raaj P et al., 2023) | 0.41 | 2 |
| AK_NLP | 0.35 | 3 |
| KEC_AI_NLP (Shanmugavadivel et al., 2023) | 0.35 | 4 |
| Athena (Sivanaiah et al., 2023a) | 0.28 | 5 |
| AbhiPaw (Bala and Krishnamurthy, 2023) | 0.27 | 6 |
| DeepBlueAI (Luo and Wang, 2023) | 0.26 | 7 |
| Habesha (Yigezu et al., 2023) | 0.22 | 8 |
| Supernova (Reddy et al., 2023) | 0.15 | 9 |

Table 5: Rank list for the Tamil-English subtask

| Team Name | macro F1 | Rank |
|---|---|---|
| DeepBlueAI (Luo and Wang, 2023) | 0.55 | 1 |
| AK-NLP | 0.51 | 2 |
| Harmony (Raaj P et al., 2023) | 0.50 | 3 |
| MUCS (Hegde et al., 2023) | 0.49 | 4 |
| Avalanche (Sivanaiah et al., 2023b) | 0.44 | 5 |
| KEC_AI_NLP (Shanmugavadivel et al., 2023) | 0.42 | 6 |
| Athena (Sivanaiah et al., 2023a) | 0.37 | 7 |
| Avalanche (Sivanaiah et al., 2023b) | 0.35 | 8 |
| CSSCUTN (Pannerselvam et al., 2023) | 0.35 | 8 |
| AbhiPaw (Bala and Krishnamurthy, 2023) | 0.29 | 9 |
| Habesha (Yigezu et al., 2023) | 0.26 | 10 |
| SuperNova (Reddy et al., 2023) | 0.25 | 11 |

Table 6: Rank list for the Telugu-English subtask

| Team Name | macro F1 | Rank |
|---|---|---|
| MUCS (Hegde et al., 2023) | 0.7459 | 1 |
| DeepBlueAI (Luo and Wang, 2023) | 0.7318 | 2 |
| Habesha (Yigezu et al., 2023) | 0.6519 | 3 |
| AK-NLP | 0.6430 | 4 |
| AbhiPaw (Bala and Krishnamurthy, 2023) | 0.6319 | 5 |
| Supernova (Reddy et al., 2023) | 0.6189 | 6 |
| Athena (Sivanaiah et al., 2023a) | 0.6137 | 7 |
| CSSCUTN (Pannerselvam et al., 2023) | 0.5939 | 8 |

Abhinaba Bala and Parameswari Krishnamurthy. 2023. Abhipaw @ abusive comment detection in tamil and telugu-dravidianlangtech. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

B Bharathi and A Agnusimmaculate Silvia. 2021. 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 and Josephine Varsha. 2022a. Ssncse nlp@ tamilnlp-acl2022: Transformer based approach for detection of abusive comment for tamil language. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 158–164.

B Bharathi and Josephine Varsha. 2022b. SSNCSE NLP@TamilNLP-ACL2022: Transformer based approach for detection of abusive comment for Tamil language. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 158–164, Dublin, Ireland. Association for Computational Linguistics.

Bharathi Raja Chakravarthi. 2023. Detection of homophobia and transphobia in Youtube comments. *International Journal of Data Science and Analytics*, pages 1–20.

Bharathi Raja Chakravarthi, Adeep Hande, Rahul Ponnusamy, Prasanna Kumar Kumaresan, and Ruba Priyadharshini. 2022a. How can we detect homophobia and transphobia? experiments in a multilingual code-mixed setting for social media governance. *International Journal of Information Management Data Insights*, 2(2):100119.

Bharathi Raja Chakravarthi, Ruba Priyadharshini, Shubanker Banerjee, Manoj Balaji Jagadeeshan, Prasanna Kumar Kumaresan, Rahul Ponnusamy, Sean Benhur, and John Philip McCrae. 2023. Detecting abusive comments at a fine-grained level in a low-resource language. *Natural Language Processing Journal*, 3:100006.

Bharathi Raja Chakravarthi, Ruba Priyadharshini, Thenmozhi Durairaj, John McCrae, Paul Buitelaar, Prasanna Kumaresan, and Rahul Ponnusamy. 2022b. Overview of the shared task on homophobia and transphobia detection in social media comments. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 369–377, Dublin, Ireland. Association for Computational Linguistics.

Vikram Gupta, Sumegh Roychowdhury, Mithun Das, Somnath Banerjee, Punyajoy Saha, Binny Mathew, Animesh Mukherjee, et al. 2022. Multilingual abusive comment detection at scale for indic languages. *Advances in Neural Information Processing Systems*, 35:26176–26191.

Asha Hegde, Sharal G, Kavya andCoelho, and Hosahalli Lakshmaiah Shashirekha. 2023. Mucs@dravidianlangtech2023: Leveraging learning models to identify abusive comments in code-mixed dravidian languages. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

Zhipeng Luo and Jiahui Wang. 2023. Deepblueai@dravidianlangtech. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

Mounika Marreddy, Subba Reddy Oota, Lakshmi Sireesha Vakada, Venkata Charan Chinni, and Radhika Mamidi. 2022. Am i a resource-poor language? data sets, embeddings, models and analysis for four different nlp tasks in telugu language. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 22(1):1–34.

Kathiravan Pannerselvam, Saranya Rajiakodi, Rahul Ponnusamy, and Sajeetha Thavareesan. 2023. Csscutn@dravidianlangtech:abusive comments 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.

SN Prasanth, R Aswin Raj, P Adhithan, B Premjith, 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.

B Premjith, KP Soman, and M Anand Kumar. 2018. A deep learning approach for malayalam morphological analysis at character level. *Procedia computer science*, 132:47–54.

Ruba Priyadharshini, Bharathi Raja Chakravarthi, Subalalitha Cn, Thenmozhi Durairaj, Malliga Subramanian, Kogilavani Shanmugavadivel, Siddhanth U Hegde, and Prasanna Kumaresan. 2022a. Overview of abusive comment detection in tamil-acl 2022. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 292–298.

Ruba Priyadharshini, Bharathi Raja Chakravarthi, Subalalitha Chinnaudayar Navaneethakrishnan, Thenmozhi Durairaj, Malliga Subramanian, Kogilavani Shanmugavadivel, Siddhanth U Hegde, and Prasanna Kumar Kumaresan. 2022b. 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*.

Amrish Raaj P, Abirami Murugappan, Lysa Packiam R S, and Deivamani M. 2023. Harmony@dravidianlangtech: Transformer-based ensemble learning for abusive comment detection. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

Ankitha Reddy, Pranav Moorthi, and Ann Maria Thomas. 2023. Supernova@dravidianlangtech 2023@abusive comment detection in tamil and telugu - (tamil, tamil-english, telugu-english). In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

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. IEEE.

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. IEEE.

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. IEEE.

Kogilavani Shanmugavadivel, Malliga Subramanian, ShriDurga R, SRIGHA S, Sree Harene J S, and Yasvanth Bala P. 2023. Kec_ai_nlp@dravidianlangtech: Abusive comment detection in tamil language. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

Rajalakshmi Sivanaiah, Angel Deborah S, M Hema, and Anza Prem. 2023a. Athena@dravidianlangtech: Abusive comment detection in code-mixed languages using machine learning techniques. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

Rajalakshmi Sivanaiah, Rajasekar S, Srilakshmisai K, Angel Deborah S, and Mirnalinee ThankaNadar. 2023b. Avalanche@dravidianlangtech:abusive comment detection in code mixed data using machine learning techniques with undersampling. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

CN Subalalitha. 2019. Information extraction framework for Kurunthogai. *Sādhanā*, 44(7):156.

Malliga Subramanian, Rahul Ponnusamy, Sean Benhur, Kogilavani Shanmugavadivel, Adhithiya Ganesan, Deepti Ravi, Gowtham Krishnan Shanmugasundaram, Ruba Priyadharshini, and Bharathi Raja Chakravarthi. 2022. Offensive language detection in Tamil youtube comments by adapters and cross-domain knowledge transfer. *Computer Speech & Language*, 76:101404.

Krithika Swaminathan, Bharathi B, Gayathri G L, and Hrishik Sampath. 2022. SSNCSE_NLP@LT-EDI-ACL2022: Homophobia/transphobia detection in multiple languages using SVM classifiers and BERT-based transformers. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 239–244, Dublin, Ireland. Association for Computational Linguistics.

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. IEEE.

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. IEEE.

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. IEEE.

Charangan Vasantharajan, Ruba Priyadharshini, Prasanna Kumar Kumarasen, Rahul Ponnusamy, Sathiyaraj Thangasamy, Sean Benhur, Thenmozhi Durairaj, Kanchana Sivanraju, Anbukkarasi Sampath, and Bharathi Raja Chakravarthi. 2022. TamilEmo: Fine-grained Emotion Detection Dataset for Tamil. In *International Conference on Speech and Language Technologies for Low-resource Languages*, pages 35–50. Springer.

Mesay Gemeda Yigezu, Selam Kanta, Grigori Kolesnikova, Olga andSidorov, and Gelbukh Alexander. 2023. Habesha@dravidianlangtech: Abusive comment detection using deep learning approach. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.