# Punny_Punctuators@DravidianLangTech-EACL2024: Transformer-based Approach for Detection and Classification of Fake News in Malayalam Social Media Text

**Nafisa Tabassum***, **Sumaiya Rahman Aodhora***, **Rowshon Akter**
**Jawad Hossain, Shawly Ahsan and Mohammed Moshiul Hoque**
Department of Computer Science and Engineering
Chittagong University of Engineering and Technology
{u1804066, u1804127, u1804003, u1704039, u1704057}@student.cuet.ac.bd
moshiul_240@cuet.ac.bd

## Abstract

The alarming rise of fake news on social media poses a significant threat to public discourse and decision-making. While automatic detection of fake news offers a promising solution, research in low-resource languages like Malayalam often falls behind due to limited data and tools. This paper presents the participation of team Punny_Punctuators in the Fake News Detection in Dravidian Languages shared task at DravidianLangTech@EACL 2024, addressing this gap. The shared task focuses on two subtasks: 1. classifying social media texts as original or fake, and 2. categorizing fake news into 5 categories. We experimented with various machine learning (ML), deep learning (DL) and transformer-based models as well as processing techniques such as transliteration. Malayalam-BERT achieved the best performance on both sub-tasks, which obtained us $2^{nd}$ place with a macro $f_1$-score of 0.87 for the subtask-1 and $11^{th}$ place with a macro $f_1$-score of 0.17 for the subtask-2. Our results highlight the potential of transformer models for low-resource languages in fake news detection and pave the way for further research in this crucial area.

## 1 Introduction

In the current digital era, fake news and the spread of false information are widespread issues that have negative effects on people, communities, and countries (Raja et al., 2023). People can be misled and deceived by fake news, which can cause them to lose trust in organizations and information sources. The dissemination of misleading narratives and the promotion of biased opinions can cause societal division and conflict (Wani et al., 2023). Therefore, using cutting-edge natural language processing (NLP) techniques, academics, policymakers, and stakeholders are working to construct strong

computational systems to prevent the dissemination of fake content (Hossain et al., 2022b).

Researchers have actively pursued the development of effective solutions to detect fake news across multiple languages in recent years (LekshmiAmmal et al., 2022). Fake news detection systems predominantly focus on high-resource languages like Spanish and English, neglecting low-resource Dravidian languages like Tulu, Malayalam, Tamil, Telugu, and Kannada due to resource scarcity (Hegde and Shashirekha, 2021). Despite the widespread use of the Malayalam language in Kerala, there is a lack of research on this language, necessitating the development of robust models for detecting fake news in Malayalam (Coelho et al., 2023). Malayalam's unique linguistic complexities, such as dialect variations, word semantics, and idiomatic expressions, make it challenging to process and analyze its text (Coelho et al., 2023). This work aims to develop a classification system for detecting fake news in Malayalam using various language technologies, aiming to identify fake articles written in the Malayalam language accurately. To effectively address the challenge, this work's major contributions are demonstrated by the following:

- Developed several machine learning (ML), deep learning (DL), and transformer-based models to identify fake news in the Malayalam language.

- Investigated and assessed the performance of the models using a variety of metrics to determine the best approach for the classification of fake news.

## 2 Related Work

The widespread use of social media and accessible internet access has resulted in the creation of mil-

---

*Authors have contributed equally to this work

lions of posts and comments per minute (Hossain et al., 2022a). Fake news on social media leads to wrong judgments, prompting studies on various machine learning and deep learning models using different word embedding techniques (Sharif et al., 2021). The best score $f_1$-score they obtained using SVM was 94.39% in task-A. Rasel et al., 2022 achieved 95.9% accuracy by building a dataset with 4678 distinct news and improved the existing dataset accuracy from 1.4% to 3.4% using CNN. Roy et al., 2019 developed CNN and BiLSTM networks individually, then fed them into a Multi-layer Perceptron Model (MLP) which resulted in 44.87% accuracy. Rai et al., 2022 proposed a BERT model with a feed-forward network with 768 hidden sizes connected to an LSTM layer for fake news classification, outperforming the vanilla pre-trained model on two datasets with a 2.50% and 1.10% increase in accuracy. Research primarily focuses on high-resource languages, with few studies on detecting abusive language in low-resource languages like Malayalam. Coelho et al., 2023 achieved 0.831 macro $f_1$-score through applying TF-IDF as a feature extraction technique and ensembling three machine learning models (MNB, LR, SVM) with majority voting. Bala and Krishnamurthy, 2023 fine-tuned the MURIL variant named "mural-base-cased" in detecting fake news in Dravidian languages resulting accuracy of 87%. Wani et al., 2023 in their work, identify toxic fake news to save time on assessing non-toxic instances. Traditional and transformer-based machine learning techniques were employed and the linear SVM method outperformed BERT SVM, RF, and BERT RF with an accuracy of 92%. Using a newly annotated dataset, another study (Chakravarthi et al., 2023) showed MURIL as a multilingual transformer by effectively detecting abusive comments in the low-resource Tamil language. Several machine learning, deep neural networks, and transformer-based approaches were utilized by Rahman et al., 2022 to analyze 5K fake news data, achieving a maximum $f_1$-score of 98% using XLM-R.

## 3   Task and Dataset Descriptions

For shared task-1, the task organizer[1] created a benchmark corpus for fake news detection (Subra-

manian et al., 2023). We utilized the corpus provided by the shared task-1 organizer to create the fake news detection classifier model. There are two subtasks in this shared task: subtask-1 and subtask-2. The goal of subtask-1 is to determine if a social media text is original or fake in the Malayalam language. In subtask-2, the objective is to create efficient models that can precisely identify and categorize fake news articles in Malayalam into five distinct classes. This subtask-2 offers five classes, including False, Half True, Mostly False, Partly False, and Mostly True. For subtask-1, the train, valid, and test splits of this dataset comprise 3257, 815 and 1019 texts respectively. Class-wise samples and dataset statistics are provided in Table 1. For Subtask 2, the train and test splits of this dataset comprise 1669 and 250 texts. Table 2 provides class-wise samples and dataset statistics of subtask-2. Because of the imbalance class distribution, we augmented the training dataset for further process.

| Classes | Train | Valid | Test | $W_T$ |
|---------|-------|-------|------|-------|
| Original | 1658 | 409 | 512 | 13268 |
| Fake | 1599 | 406 | 507 | 21420 |
| **Total** | **3257** | **815** | **1019** | **34688** |

Table 1: Class-wise distribution of train, validation and test set for subtask-1, where $W_T$ denotes total words in the Train dataset

| Classes | Train | Test | $W_T$ |
|---------|-------|------|-------|
| False | 1246 | 149 | 12183 |
| Mostly False | 239 | 63 | 2380 |
| Half True | 141 | 24 | 1462 |
| Partly False | 42 | 14 | 363 |
| Mostly True | 1 | 0 | 8 |
| **Total** | **1669** | **250** | **16396** |

Table 2: Class-wise distribution of train and test set for subtask-2, where $W_T$ denotes total words in Train dataset

In the pre-processing stage, we filter out URLs, emojis, and assorted symbols present within the provided dataset.

## 4   Methodology

This section provides a brief overview of the methods and approaches used to solve the problem mentioned in the previous section. Several machine learning and deep learning approaches were em-

---

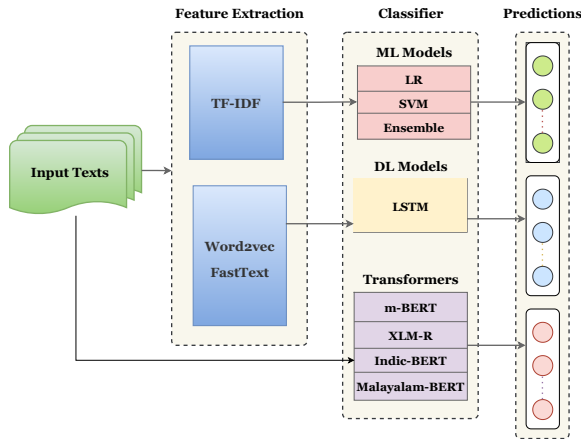[1]https://sites.google.com/view/dravidianlangtech-2024/shared-tasks-2024?authuser=0

Figure 1: Abstract process of fake news detection

ployed for developing the baseline models. Performance of the classification of both tasks improved significantly by fine-tuning the transformer-based models. To identify fake news, this work employed four pre-trained transformer-based models: XLM-R, Malayalam BERT, mBERT, and Indic-BERT. To be more specific, we fine-tuned the Huggingface transformer[2] library's "XLM-RoBERTa-base" (Conneau et al., 2019), "l3cube-pune / Malayalam BERT"(Joshi, 2022), "BERT-base-multilingual-cased"(Devlin et al., 2018) and "Indic-BERT"(Kakwani et al., 2020). Figure 1 depicts the developed system's schematic process.

### 4.1 Pre-processing

During pre-processing, noises like punctuation, alphanumeric letters, and special characters (slash, brackets, ampersands, etc.) are eliminated from Malayalam code-mixed data and transliterated (Raihan et al., 2023).

### 4.2 Machine Learning Models

The feature vector has been extracted for machine learning techniques using Word2Vec and TF-IDF word embedding in subtask-1 (Mikolov et al., 2013). To establish a fake news detection system, we start by using conventional machine learning techniques like Logistic Regression (LR), Support Vector Machine (SVM) and Random Forest (RF). The application of an ensemble of machine learning classifiers was further enhanced for improved results. Ensemble technique was explored by subtask-1 using Random Forest (RF) and Decision Tree (DT) classifiers in conjunction with LR and SVM,

[2]https://huggingface.co/docs/transformers/index

whereas 'n estimators = 100' were utilized for DT and RF. Majority voting technique is applied to achieve the prediction from the ensemble method.

### 4.3 Deep Learning Models

For classification tasks, DL algorithms performed better. Two different deep learning models were used to classify fake news: LSTM (Word2Vec)and LSTM (FastText) in subtask-1. LSTM is widely recognized for its proficiency in capturing both the semantic details and sustained dependencies over the long term. While Bidirectional LSTM (BiLSTM) takes advantage of both past and future states, LSTM records semantic information and long-term dependencies. The models were trained using the Adam optimizer with a learning rate of 1e-3 and batch size of 32. To obtain the predictions, a sigmoid layer was employed at the end.

### 4.4 Transformer Models

In the past few years, transformers have gained increasing recognition due to their exceptional capabilities across different areas of natural language processing (NLP). We investigated four transformers XLM-R, m-BERT, Malayalam-BERT, and Indic-BERT pre-trained models to fine-tune on Malayalam fake news detection dataset. A self-supervised training method for cross-lingual comprehension, known as XLM-R (Conneau et al., 2019), is especially useful for low-resource languages. A transformer model called m-BERT (Devlin et al., 2018) has been pre-trained on 104 different languages; Malayalam-BERT (Joshi, 2022) has been pre-trained only on the Malayalam language. A multilingual ALBERT model embracing 12 main Indian languages, IndicBERT (Kakwani et al., 2020) were trained on a large corpus. The ktrain package is used to refine these models, which are selected from the Pytorch Huggingface transformers library. Each model has been trained for a maximum of 10 epochs using a batch size of 9 for Malayalam-BERT and 16 for XLM-R and m-BERT.

## 5 Results and Analysis

This section evaluates the performance of different models in recognizing fake news detection in Dravidian languages.

### 5.1 Evaluation Metrics

Model performance was measured using the macro-averaged F1 score. Table 3 and Table 4 present the

| Classifiers | Original | | | Transliterated | | |
|---|---|---|---|---|---|---|
| | P | R | F | P | R | F |
| LR | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 |
| SVM | 0.80 | 0.80 | 0.80 | 0.83 | 0.83 | 0.83 |
| RF | 0.77 | 0.77 | 0.77 | 0.83 | 0.83 | 0.83 |
| Ensemble | 0.79 | 0.80 | 0.78 | 0.83 | 0.83 | 0.83 |
| LSTM (Word2vec) | 0.80 | 0.79 | 0.79 | 0.84 | 0.85 | 0.84 |
| LSTM (Fasttext) | 0.78 | 0.79 | 0.78 | 0.81 | 0.81 | 0.81 |
| Indic-BERT | 0.75 | 0.75 | 0.75 | 0.81 | 0.81 | 0.81 |
| M-BERT | 0.84 | 0.84 | 0.84 | 0.85 | 0.85 | 0.85 |
| XLM-R | 0.86 | 0.86 | 0.86 | 0.87 | 0.87 | 0.87 |
| **Malayalam-BERT** | **0.86** | **0.86** | **0.86** | **0.87** | **0.87** | **0.87** |

Table 3: Performance of various models on the test set of subtask-1. Here P, R, and F1 denotes macro Precision, macro Recall, and macro F1-Score respectively.

| Classifiers | Original | | | Transliterated | | |
|---|---|---|---|---|---|---|
| | P | R | F | P | R | F |
| Indic-BERT | 0.14 | 0.25 | 0.12 | 0.02 | 0.25 | 0.04 |
| XLM-R | 0.14 | 0.25 | 0.12 | 0.07 | 0.25 | 0.11 |
| M-BERT | 0.12 | 0.25 | 0.14 | 0.21 | 0.23 | 0.16 |
| **Malayalam-BERT** | **0.22** | **0.20** | **0.16** | **0.17** | **0.21** | **0.17** |

Table 4: Performance of various models on the test set of subtask-2. Here P, R, and F1 denotes macro Precision, macro Recall, and macro F1-Score respectively.

performance of each classifier, trained on both the original and augmented training datasets.

## 5.2 Comparative Analysis

Regarding subtask-1, the $f_1$-scores highlight LR model sustained proficiency with a consistent $f_1$-score of 0.82 in both original and transliterated datasets. In the original dataset, SVM and RF demonstrate parallel performance, achieving $f_1$-scores of 0.80 and 0.77, respectively. Interestingly, in the transliterated dataset, both SVM and RF experience a slight boost in $f_1$-score to 0.83. The ensemble method exhibits robustness across both datasets, maintaining $f_1$-score of 0.79 in the original dataset and a commendable improvement (0.83) in the transliterated dataset.

Moving to deep learning techniques, LSTM with Word2Vec consistently outperformed others, achieving macro $f_1$-scores of 0.79 and 0.84 on the original and transliterated sets, respectively. Fast-Text, however, yielded lower scores of 0.78 and 0.81 in the original and augmented datasets.

Transformer models outperformed both machine learning and deep learning models in each task. On the original dataset, all transformers, except Indic-BERT, surpassed the highest macro $f_1$-score (0.82)

achieved by SVM. Malayalam-BERT emerged as the top performer with a leading score of 0.86. In the transliterated set, excluding Indic-BERT, all transformers excelled beyond the 0.82 macro $f_1$-score from machine learning and deep learning models, with Malayalam-BERT achieving the highest macro $f_1$-score of 0.87 in the augmented training set. This suggests the consistent superiority of transformers, particularly Malayalam-BERT, in Malayalam fake news identification.

In subtask-2, with the original dataset, Indic-BERT and XLM-R exhibit similar macro $f_1$-score of 0.12, indicating relatively low performance overall, while M-BERT and Malayalam-BERT achieved macro $f_1$-scores of 0.14 and 0.16, respectively. However, in the augmented dataset, there are notable improvements for some classifiers. M-BERT and Malayalam-BERT show increased macro $f_1$-scores of 0.16 and 0.17. respectively, suggesting a boost in performance. Malayalam-BERT remains the best performer even in the transliterated dataset.

## 5.3 Error Analysis

A thorough examination of error analysis is conducted both quantitatively and qualitatively to offer

Figure 2: Confusion matrix of the best model, Malayalam-BERT, after running on the subtask-1 dataset
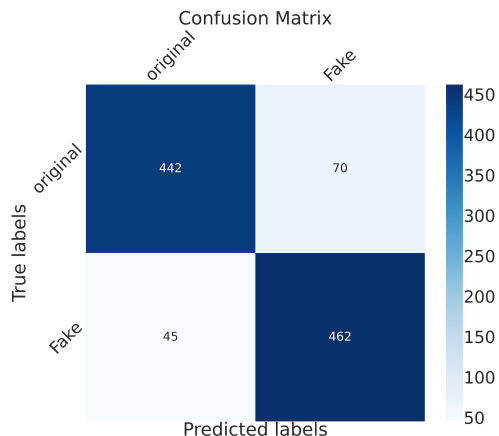


Figure 3: Confusion matrix of the best model, Malayalam-BERT, after running on the subtask-2 dataset

comprehensive insights into the effectiveness of the suggested model.

### 5.3.1 Quantitative Analysis

We can see the best performing model was Malayalam-BERT in subtask-1 which produced an $f_1$ score of 0.87 from Table 3. Figure 2 represents the confusion matrix of this model. From the confusion matrix, we can carry out an error analysis for this model. Among the two classes, the fake class has the highest True-Positive Rate (TPR) of 91.12%, while the original class has a lower TPR of 86.32%.

Figure 3 represents the confusion matrix of subtask-2. The FALSE class dominates with the highest True Positive Rate (TPR), closely followed by HALF TRUE. However, the discrimination is glaring as MOSTLY FALSE and MOSTLY TRUE struggle to detect any data accurately. Particularly, PARTLY FALSE suffers from a remarkably low detection rate for the True class. This bias stems from the dataset's imbalance, accentuating the need for a more equitable distribution for improved model performance.

### 5.3.2 Qualitative Analysis

Figure 4 and 5 in Appedix 8 provide illustrations of predicted outcomes by the best-performing model, Malayalam-BERT.
The model accurately predicts the outcomes for the 1st, 2nd, and 5th text samples in subtask-1. However, it encounters challenges in predicting the 3rd and 4th text samples, where its performance is not as successful for subtask-1. In subtask-2 the model correctly predicts text samples 1 and 2,
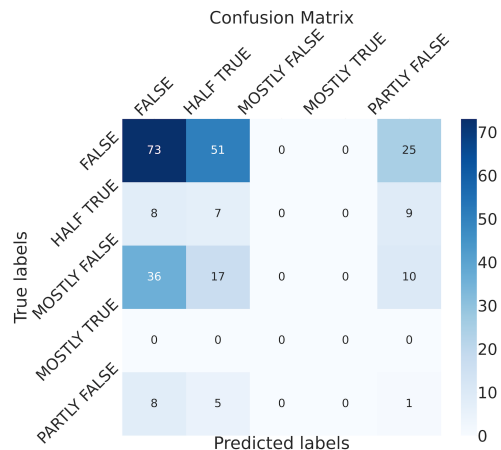
however, it has trouble predicting text samples 3 and 4. The problem of class imbalance might be the cause of incorrect predictions.

## 6 Limitations

Our models depend on accurately turning words from one language into another (transliteration). Even small mistakes can significantly impact the information, affecting how well our models work. Furthermore, in subtask-2, the data was quite imbalanced, and no methods were applied to balance the data. Future exploration of methods with multiple transliteration options, as well as augmentation for balancing the dataset can be investigated to further enhance the accuracy of our approach.

## 7 Conclusion

This work aimed to detect and classify fake news from Malayalam social media text. We have thoroughly investigated several machine learning (ML), and deep learning (DL) and transformer-based models for Malayalam fake news identification and classification. The Malayalam-BERT model has proven to be more effective than the others, as evidenced by its highest macro F1-Score of 0.87 for subtask-1 and 0.17 for subtask-2. In subtask-1, the model excels, securing the 2nd position with a noteworthy macro $f_1$-score of 0.87. In contrast, in subtask-2, it ranks 11th with a macro $f_1$-score of 0.17. To improve model performance in the future, we want to look at different architectures and use ensemble techniques. We'll explore different ways to tackle problems that come from imbalanced datasets.

# References

Abhinaba Bala and Parameswari Krishnamurthy. 2023. Abhipaw@ dravidianlangtech: Fake news detection in dravidian languages using multilingual bert. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, pages 235–238.

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.

Sharal Coelho, Asha Hegde, G Kavya, and Hosahalli Lakshmaiah Shashirekha. 2023. Mucs@ dravidianlangtech2023: Malayalam fake news detection using machine learning approach. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, pages 288–292.

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. *arXiv preprint arXiv:1911.02116*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Asha Hegde and Hosahalli Lakshmaiah Shashirekha. 2021. Urdu fake news detection using ensemble of machine learning models. In *CEUR Workshop Proceedings*, pages 132–141.

Alamgir Hossain, Mahathir Bishal, Eftekhar Hossain, Omar Sharif, and Mohammed Moshiul Hoque. 2022a. Combatant@ tamilnlp-acl2022: Fine-grained categorization of abusive comments using logistic regression. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 221–228.

Alamgir Hossain, Mahathir Bishal, Eftekhar Hossain, Omar Sharif, and Mohammed Moshiul Hoque. 2022b. COMBATANT@TamilNLP-ACL2022: Fine-grained categorization of abusive comments using logistic regression. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 221–228, Dublin, Ireland. Association for Computational Linguistics.

Raviraj Joshi. 2022. L3Cube-MahaCorpus and MahaBERT: Marathi monolingual corpus, Marathi BERT language models, and resources. In *Proceedings of the WILDRE-6 Workshop within the 13th Language Resources and Evaluation Conference*, pages 97–101, Marseille, France. European Language Resources Association.

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, NC Gokul, Avik Bhattacharyya, Mitesh M Khapra, and Pratyush Kumar. 2020. Indicnlpsuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for indian languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4948–4961.

Hariharan LekshmiAmmal, Manikandan Ravikiran, and Anand Kumar Madasamy. 2022. NITK-IT_NLP@TamilNLP-ACL2022: Transformer based model for toxic span identification in Tamil. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 75–78, Dublin, Ireland. Association for Computational Linguistics.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26.

MD Sijanur Rahman, Omar Sharif, Avishek Das, Sadia Afroze, and Mohammed Moshiul Hoque. 2022. Fand-x: Fake news detection using transformer-based multilingual masked language model. In *2022 IEEE International Women in Engineering (WIE) Conference on Electrical and Computer Engineering (WIECON-ECE)*, pages 153–158. IEEE.

Nishant Rai, Deepika Kumar, Naman Kaushik, Chandan Raj, and Ahad Ali. 2022. Fake news classification using transformer based enhanced lstm and bert. *International Journal of Cognitive Computing in Engineering*, 3:98–105.

Md Nishat Raihan, Umma Hani Tanmoy, Anika Binte Islam, Kai North, Tharindu Ranasinghe, Antonios Anastasopoulos, and Marcos Zampieri. 2023. Offensive language identification in transliterated and code-mixed bangla. *arXiv preprint arXiv:2311.15023*.

Eduri Raja, Badal Soni, and Sami Kumar Borgohain. 2023. nlpt malayalm@ dravidianlangtech: Fake news detection in malayalam using optimized xlm-roberta model. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, pages 186–191.

Risul Islam Rasel, Anower Hossen Zihad, Nasrin Sultana, and Mohammed Moshiul Hoque. 2022. Bangla fake news detection using machine learning, deep learning and transformer models. In *2022 25th International Conference on Computer and Information Technology (ICCIT)*, pages 959–964. IEEE.

Arjun Roy, Kingshuk Basak, Asif Ekbal, and Pushpak Bhattacharyya. 2019. A deep ensemble framework for fake news detection and multi-class classification of short political statements. In *Proceedings of the 16th International Conference on Natural Language Processing*, pages 9–17.

Omar Sharif, Eftekhar Hossain, and Mohammed Moshiul Hoque. 2021. Combating hostility: Covid-19 fake news and hostile post detection in social media. *CoRR*, abs/2101.03291.

Malliga Subramanian, Bharathi Raja Chakravarthi, Kogilavani Shanmugavadivel, Santhiya Pandiyan, Prasanna Kumar Kumaresan, Balasubramanian Palani, Muskaan Singh, Sandhiya Raja, Vanaja, and Mithunajha S. 2023. Overview of the shared task on fake news detection from social media text. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*, Varna, Bulgaria. Recent Advances in Natural Language Processing.

Mudasir Ahmad Wani, Mohammad ELAffendi, Kashish Ara Shakil, Ibrahem Mohammed Abuhaimed, Anand Nayyar, Amir Hussain, and Ahmed A Abd El-Latif. 2023. Toxic fake news detection and classification for combating covid-19 misinformation. *IEEE Transactions on Computational Social Systems*.

# 8 Appendix

| Text Sample | Actual | Predicted |
|---|---|---|
| മൊഴി മുത്തുകളായ വരികൾ 😂😂(Lines that are pearls of expression) | Original | Original |
| Shame for entire Woman&#39 | Original | Original |
| വിമർശിക്കുന്നു(Criticizing) | Fake | Original |
| They revised the social distancing from 2M to 2cm so it's fine. 🔥😡 | Original | Fake |
| എല്ലാം വരുടെയും (പ്രാത്ഥന മൂലം ആണ്) (Everything is due to the prayer of the bridegroom) | Fake | Fake |

Figure 4: Few examples of predicted outputs by the proposed (Malayalam-BERT) model for subtask-1

| Text Sample | Actual | Predicted |
|---|---|---|
| കീർത്തി സുരേഷ് ഫർഹാൻ എന്ന മുസ്ലിം യുവാവിനെ കല്യാണം കഴിയ്ക്കുന്നു (Keerthi is getting married to a young Muslim named Suresh Farhan) | False | False |
| വാരിയം കുന്നനെ അറസ്റ്റ് ചെയ്തതായി മലയാള മനോരമ നൽകിയ വാർത്ത (Malayalam Manorama reported that Variyam Kunnan was arrested.) | Mostly False | Mostly False |
| ചന്ദനക്കുറിയണിഞ്ഞ് വിഎസ് അച്ചുതാനന്ദൻ.(VS Achuthanandan dressed in sandalwood.) | False | Mostly True |
| പലസ്തീൻ പതാകയണിഞ്ഞ് ക്രിസ്ത്യാനോ റൊണാൾഡോ (Cristiano Ronaldo wearing the Palestinian flag) | Half True | Mostly False |

Figure 5: Few samples of the predicted outcomes of the proposed (Malayalam-BERT) model for subtask-2