

# A Novel Explainable AI Framework for Detecting Misinformation on Federated Social Media Data

P.Lokanadham<sup>1</sup>, Punuru Kavyashree<sup>2</sup>, S.Khuranun Mubeen<sup>3</sup>, N.Venu Madhav<sup>4</sup>, K.Deepak Raj<sup>5</sup>

<sup>1</sup>Assistant Professor, Dept of Information Technology, SV College of Engineering, Tirupati, India

<sup>2</sup>B. Tech, Dept of Information Technology, SV college of Engineering, Tirupati, India.

<sup>3</sup>B. Tech, Dept of Information Technology, SV college of Engineering, Tirupati, India.

<sup>4</sup>B. Tech, Dept of Information Technology, SV college of Engineering, Tirupati, India.

<sup>5</sup>B. Tech, Dept of Information Technology, SV college of Engineering, Tirupati, India.

## Abstract

Web Information Processing (W.I.P.) has enormously impacted modern society since a huge percentage of the population relies on the internet to acquire information. Social Media platforms provide a channel for disseminating information and a breeding ground for spreading misinformation, creating confusion and fear among the population. One of the techniques for the detection of misinformation is machine learning-based models. However, due to the availability of multiple social media platforms, developing and training AI-based models has become a tedious job. Despite multiple efforts to develop machine learning-based methods for identifying misinformation, more work must be done on developing an explainable generalized detector capable of robust detection and generating explanations beyond black-box outcomes. Knowing the reasoning behind the outcomes is essential to make the detector trustworthy. Hence employing explainable A.I. techniques is of utmost importance.

In this work, the integration of two machine learning approaches, namely domain adaptation and explainable A.I., is proposed to address these two issues of generalized detection and explainability. Firstly, the Domain Adversarial Neural Network (DANN) develops a generalized misinformation detector across multiple social media platforms. DANN generates the classification results for test domains with relevant but unseen data. The DANN-based traditional black-box model cannot justify and explain its outcome, i.e., the labels for the target domain. Hence a Local Interpretable Model-Agnostic Explanations (LIME) explainable A.I. model is applied to explain the outcome of the DANN model. To demonstrate these two approaches and their integration for effective explainable generalized detection, COVID-19 misinformation is considered as a case study. We experimented with two datasets and compared results with and without DANN implementation. It is observed that using DANN significantly improves the F1 score of classification and increases the accuracy by 3% and AUC by 9%. The results show that the proposed framework outperforms well in the case of domain shift and can learn domain-invariant features while explaining the target labels with LIME implementation. This can enable trustworthy information processing and extraction to combat misinformation effectively.

**Keywords:** Covid19, DANN, LIME, Misinformation Detection, Social Media, Text Processing, Web Information Processing, XAI.