# IndDeepFake: Mitigating the Spread of Misinformation in India through a Multimodal Adversarial Network

**Manish Kumar Singh[1], Jawed Ahmed[1], Kamlesh Kumar Raghuvanshi[2], Mohammad Afshar Alam[1]**

**Abstract:** The issue of the distribution of fake news and misinformation on social media platforms is a rising global concern, which has affected India as well. This research paper introduces an innovative method for identifying and detecting fake news in India using a multimodal adversarial network. The approach presented in this study leverages both text and image characteristics to encompass the multimodal aspects of fake news. Adversarial training is employed to learn robust and discriminative features/characteristics that enable differentiating authentic news from fabricated news. Evaluation of the proposed method is conducted on an Indian fake news events dataset and achieves a high accuracy and  F1-score of 0.89 and 0.90 respectively. The experiment results indicate that the proposed multimodal adversarial network approach is effective in detecting fake news in the Indian context and thus helpful in mitigating the dissemination of misinformation.

*Keywords: Fake news, Multimodal adversarial network, Misinformation mitigation, Indian fake news dataset, Information Credibility*

## 1.    Introduction

The issue of dissemination of fake news is growing concern in numerous nations, including India, where it can have serious consequences for society. The possible definition of fake news is the "fabricated information that mimics news media content in form but not in organizational process or intent" [1]. It can spread quickly through social media & other channels and can be challenging to differentiate it from genuine news. Fake news can have a variety of negative impacts, including inciting violence, damaging reputations, and influencing political decisions [2].

Various methods have been proposed for identifying false information, such as machine learning algorithms, natural language processing techniques, and network analysis. These methods have been developed based on different features such as linguistic, behavioral, and contextual characteristics of fake news [3].

A commonly employed method involves utilizing machine learning algorithms, like support vector machines (SVMs) and deep neural networks (DNNs), to categorize disseminated news as authentic or false by analyzing their content. Researchers have experimented with different feature sets, including textual features, social network features, and user engagement features, to enhance the accuracy of these techniques [4].

Another approach is to use network analysis techniques to identify fake news based on its propagation pattern on social media. These methods leverage network properties such as centrality and community structure to detect fake news that has been spread by bots or coordinated groups of users [5].

However, despite their potential, these methods often have limitations. For example, machine learning algorithms may have difficulty in detecting sophisticated fake news that is designed to evade detection by mimicking the style and structure of real news articles [6]. Similarly, network analysis techniques may be limited by the reliability and quantity of available data, as well as the dynamic nature of social networks [7].

To overcome these limitations, an adversarial multimodal network for detecting false news is presented in the current research paper. The network so proposed draws inspiration from the achievements of adversarial networks in different fields [8], such as generative models and image synthesis. The network of the proposed method consists of a (a). generator network, responsible for producing fabricated multimodal data, (b). a discriminator network, which distinguishes between genuine and fabricated data, and (c). a classifier network, which determines whether the data input is real or constitutes fake news. Through training the network on a dataset consisting of both real and fabricated multimodal data, this research demonstrates that the proposed technique achieves superior accuracy in identifying fake news when compared to existing methods.

This article presents an innovative method to identify false information in India using a multimodal adversarial network. The proposed approach makes the following

[1]*Department of Computer Science, Jamia Hamdard University, Delhi, India.*

[2]*Department of Computer Science, Ramanujan College, University of Delhi, Delhi, India.*

[1]*Corresponding author: manish.kumar2191989@gmail.com*

contributions:

- The incorporation of text and image attributes to capture the multimodal characteristics of fabricated information specific to the Indian context.

- The use of adversarial training to learn robust and discriminative features that can differentiate between genuine and fabricated news.

- The application of the proposed approach to an Indian fake news events dataset, which has not been explored in prior studies on identifying fake news.

- The assessment of the recommended method utilizes commonly employed performance metrics such as accuracy, precision, recall, and F1 score.

The rest of the article is structured as follows: Section 2 presents a detailed review of related work that addresses fake news and misinformation. The methodology adopted for this research is explained in Section 3 and gives comprehensive information on the proposed multimodal adversarial network's architecture.. The experimental setup details, including the Indian fake news events dataset and the preprocessing steps applied to it, are presented in Section 4. Section 5 provides a discussion of the results, including performance comparison with baseline techniques, performance comparison using hyperparameters, performance comparison during the ablation study, performance comparison with sophisticated techniques for identifying false news, and sensitivity analysis of the proposed model, while pointing out the drawbacks of the suggested strategy. The article is concluded in Section 6 along with some ideas for further research.

## 2. Related Work

### 2.1 Natural Language Processing (NLP)

In recent years, there has been a lot of study on the use of NLP approaches for detecting fake news. Ma et al. (2018) [9] used sentiment analysis and clustering to detect false news propagated over social media. Gangireddy et al. (2020) [10] presented a technique based on graphs that leverages user behavior and content characteristics for detecting false news. Raza and Ding (2022) [11] developed a transformer-based model to identify false news by utilizing its headline and body text. However, NLP techniques have limitations, especially when dealing with language ambiguity and sarcasm [12].

### 2.2 Machine Learning (ML)

ML techniques are quite popular for spotting false news. Shu et al. (2017) [3] used ML techniques to detect fake news and propaganda propagated over social media. Singhal et al. (2022) [13] proposed a model based on deep learning that uses linguistic and visual cues to detect false

news. Sahoo et al. (2021) [14] developed a deep learning-based approach that incorporates external knowledge sources in order to increase the effectiveness of false news identification. Althobaiti (2022) [15] presents a BERT-based method that makes use of emojis and emotion analysis to spot hate speech and objectionable words in Arabic tweets. The study gives a thorough methodology and results based on experiments to prove the efficacy of the suggested strategy. However, when faced with new data, ML models' performance may suffer since they need a lot of labeled data to train on [16].

### 2.3 Network Analysis

Network analysis has been another popular approach for detecting fake news. This approach leverages interpersonal relationships between users and the propagation patterns of news articles to identify fake news. Vosoughi et al. (2018) [4] analyzed networks to find the propagation of misleading news on Twitter. Nasir et al. (2021) [17] suggested a technique to identify false news that merges network analysis and deep learning. Choudhary and Arora (2021) [18] developed a model that leverages social and topical features to spot false news disseminated over social media. However, network analysis requires access to the social connections between users, which may not always be available [19].

### 2.4 Adversarial Network

Adversarial network-based approaches have been used for detecting fake news as they can learn to distinguish between real and fake news articles using a small amount of labeled data. Wang et al. (2019) [20] proposed a GAN-based model for fake news detection. Peng and Xintong (2022) [21] developed an adversarial learning-based model that uses both textual and visual features for fake news detection. Wei et al. (2022) [22] proposed a GAN-based model that leverages text, image, and metadata information to detect fake news. However, adversarial network-based models may suffer from the problem of adversarial examples, where small perturbations to the input can cause the model to misclassify the news article [23].

### 2.5 Multimodal Adversarial Network

Multimodal adversarial networks have shown promising results in detecting fake news by leveraging both textual and visual cues. Khattar et al. (2019) [24] developed a multimodal adversarial learning-based model that uses both textual and visual features for fake news detection. Quan et al. (2021) [25] proposed a multimodal GAN-based model that exploits visual, textual, and metadata features for fake news detection. Yuan et al. (2021) [26] proposed a multimodal approach based on graph neural networks and GANs for fake news detection. This approach achieved state-of-the-art performance on

multiple benchmark datasets. However, the availability of labeled multimodal datasets for training is limited. Furthermore, adversarial attacks can also be applied to multimodal data, which may affect the model's performance [23].

This article presents an innovative technique based on a multimodal adversarial network, IndDeepFake, for detecting Indian fake news events, which exploits both textual and visual cues. The proposed approach overcomes the limitations of existing techniques by effectively detecting adversarial attacks on multimodal data.

## 3. Methodology

### 3.1 Problem Statement

Given a set of news articles $N = \{n_1, n_2, ..., n_m\}$ and their corresponding labels $Y = \{y_1, y_2, ..., y_m\}$, where $y_i$ is either 0 or 1 and denotes if the news item $n_i$ is authentic or not, the proposed work aims to grasp a function f(n) that maps an input news article 'n' to a predicted label $\hat{y}$, $\hat{y} = f(n)$. The suggested technique seeks to make use of the textual as well as visual components of news items to enhance detection performance. Specifically, a multimodal adversarial network is used in the current paper to jointly learn representations from textual and visual modalities and detect adversarial attacks on multimodal data.

### 3.2 Proposed Model

The proposed multimodal adversarial network, IndDeepFake, comprises key modules including (a). a text encoder, (b). an image encoder, and (c). a classifier. The visual and linguistic elements of the news stories are encoded by the image encoder and text encoder respectively, while the classifier, based on the encoded attributes, predicts the label of the news article.

#### 3.2.1. Text Encoder

The text encoder is a deep neural network that maps a news article $x_i$ to a fixed-dimensional vector $z_i$ in $R^d$. A pre-trained transformer-based model, BERT, encodes the news article.

$$z_i = E(x_i) \qquad (1)$$

where 'E' is the transformer-based model that has already been trained.

#### 3.2.2. Image Encoder

Deep convolutional neural networks serve as the image encoder. that maps an image $I_i$ to a fixed-dimensional vector $v_i$ in $R^d$. A pre-trained CNN is used to encode the image.

$$v_i = F(I_i) \qquad (2)$$

where F is the pre-trained CNN.

#### 3.2.3. Classifier

The classifier, which is a feedforward neural network, generates a probability score for each class after receiving the concatenated text and picture characteristics as input. For the current work, a sigmoid activation function at the output layer is used to ensure that the output is between 0 and 1.

$$y_i = \sigma \ (W[g(z_i,v_i)]+b) \qquad (3)$$

where 'W' and 'b' stand for the classifier's weight and bias parameters, 'g' is a function that concatenates the text and image features, and 'σ' is the sigmoid activation function.

#### 3.2.4. Adversarial Network

To overcome the limitations of existing techniques, an adversarial network is introduced in the proposed technique that generates adversarial examples to fool the classifier. The adversarial network consists of a text adversarial network and an image adversarial network, both of which are trained to generate perturbations to the input that maximally increases the loss of the classifier.

The text adversarial network receives the encoded text feature $z_i$ as input and generates a perturbation vector $\hat{z}$ in $R^d$ that is added to the encoded text feature to generate the adversarial example:

$$\hat{z}_i = z_i + \delta_z \qquad (4)$$

The text adversarial network is trained to maximize the loss of the classifier on the generated adversarial example.

The image adversarial network receives the encoded image feature $v_i$ as input and generates a perturbation vector $\tilde{v}$ in $R_d$ that is added to the encoded image feature to generate the adversarial example.

$$\tilde{v}_i = v_i + \delta_v \qquad (5)$$

The image adversarial network is trained to maximize the loss of the classifier on the generated adversarial example.

#### 3.2.5. Training Objective

The overall training objective is to minimize classification loss while maximizing adversarial loss. In the current technique, a weighted sum of the adversarial loss and the binary cross-entropy loss is used as the training objective.

The binary cross-entropy loss calculates the discrepancy between the real label of the news article and the anticipated probability score. In order to appropriately categorize the news stories, the classifier is trained using this loss.. Let $y_i$ denote the true label for the news article $x_i$, and $p_i$ be the predicted probability of the article being fake, where $p_i = f(\theta, x_i)$ and $f\{\theta\}$ is the function that maps an input article to a probability score. Following is the definition of the binary cross-entropy loss function [8]:

$$L_{CE} = -\frac{1}{n}\sum_{i=1}^{n}(y_i * \log p_i + (1 - y_i) * \log (1 - p_i))$$
(6)(i)

The adversarial loss measures the difference between the probability scores of the classifier on the original and adversarial example. The text adversarial network is trained using this loss and the image adversarial network generates adversarial examples in order to fool the classifier. Let $\ddot{x}_i$ be the adversarial example generated by the adversarial network for the news article $x_i$, and let $\acute{p}_i = f(\theta, \ddot{x}_i^{*})$ be the predicted probability of the adversarial example being fake. Following is the definition of the adversarial loss [8]:

$$L_{adv} = -\frac{1}{n}\sum_{i=1}^{n}(y_i * \log \acute{p}i + (1 - y_i) * \log (1 - \acute{p}i))$$
(6)(ii)

The overall training objective is computed as a weighted sum of the binary cross-entropy loss and the adversarial loss which can be represented by the equation:

$$L = \alpha L_{CE} + \beta L_{adv} \qquad (7)$$

Here, $L_{CE}$ is the classification loss, $L_{adv}$ is the adversarial loss, and $\alpha$ and $\beta$ are hyperparameters that regulate the two loss terms' respective weights. The binary cross-entropy difference between the anticipated label $y_i$ and the actual label $y_i^{*}$ constitutes the classification loss. The adversarial loss is the difference between the classification score for the adversarial example and the original example, averaging over all examples in the training set.

### 3.3 *Architecture*

The architectural layout of the IndDeepFake model that is suggested in the article is shown in the flowchart of Fig. 1. The flowchart begins with a news item (X) and a label (Y) designating whether it is authentic or not. The news article is then encoded into feature vectors z and v using a text encoder E to generate, and an image encoder F for the textual and visual parts, respectively. A joint representation (g(z,v)) of the text and image features is created by concatenating the features of text and images. This representation is then used by the classifier C to predict the label of the news article.



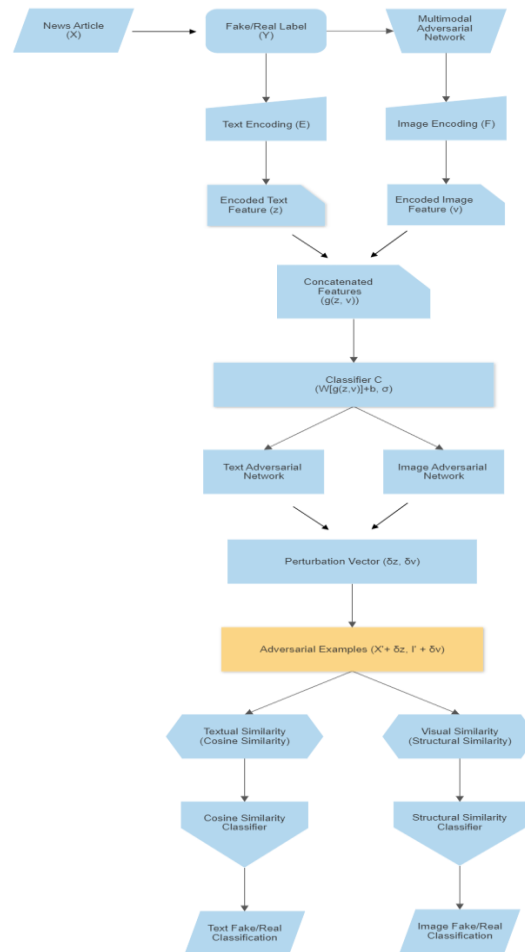**Fig. 1** Proposed Architecture of IndDeepFake

Additionally, the joint representation is used by the two separate adversarial networks, one for text and one for images. The adversarial networks are used to generate

perturbations ($\delta z$, 0) and (0, $\delta v$) to the text and image features respectively that are added to the original features 'z' and 'v' to create adversarial examples (X'+ $\delta z$, I') and

(X', I' + δv), respectively. The adversarial examples are used to train the classifier to be more robust to adversarial attacks.

The flowchart also includes two similarity classifiers, one for textual similarity (using cosine similarity) and one for visual similarity (using structural similarity). The similarity classifiers are used to compare the original and adversarial examples to identify the degree to which they resemble one another. During the final phase, the output of the similarity classifiers and the original & adversarial examples are utilized to determine if a news story is authentic or not using the loss function of the binary-cross entropy. The overall training objective is to minimize the sum of the adversarial loss and the binary cross-entropy loss, which is achieved through backpropagation and gradient descent.

### 3.4 *Mathematical Model*

Let X be the news article, Y be the binary label indicating if the article is real or fake, E be the text encoder, F be the image encoder, G be the concatenated feature generator, W be the classifier weights, b be the classifier bias, σ be the sigmoid activation function, $D_z$ be the text adversarial network, and $D_v$ be the image adversarial network.

Let z be the encoded text feature vector, v be the encoded image feature vector, $\delta_z$ be the perturbation vector for the text adversarial network, and $\delta_v$ be the perturbation vector for the image adversarial network. Let X' and I' be the adversarial examples generated by the text and image adversarial networks, respectively.

The key training objective is to reduce the subsequent loss function:

$$L = L_{CE}(Y, \sigma(W[G(z,v)]+b)) + \lambda_1 L_{adv}(D_z(z), z + \delta_z) + \lambda_2 L_{adv}(D_v(v), v + \delta_v) + \lambda_3 L_{sim}(G(z,v), G(z+\delta_z,v+\delta_v)) \quad (8)$$

where:

● $L_{CE}$ is binary cross-entropy loss which is basically the difference between the classifier output (W[G(z,v)]+b) and the real/fake label Y;

● $L_{adv}$ is the adversarial loss for the text and image adversarial networks, which is the $L_{CE}$ loss between the predicted adversarial label and the actual label (1 for text adversarial network and 0 for image adversarial network);

● $L_{sim}$ is the similarity loss in the feature space produced by G between the actual and adversarial examples; and

● $\lambda_1$, $\lambda_2$, and $\lambda_3$ are hyperparameters controlling the relative importance of each loss.

The adversarial loss for the text and image adversarial networks can be written as [27]:

$$L_{adv}(D_z(z), z + \delta_z) = -Y \log(D_z(z)) - (1-Y) \log(1 - D_z(z + \delta_z)) \quad (9)(i)$$

$$L_{adv}(D_v(v), v + \delta_v) = -Y \log(D_v(v)) - (1-Y) \log(1 - D_v(v + \delta_v)) \quad (9)(ii)$$

The similarity loss can be written as:

$$L_{sim}(G(z,v), G(z+\delta_z,v+\delta_v)) = \|G(z,v) - G(z+\delta_z,v+\delta_v)\|^2 \quad (10)$$

The overall objective is optimized using gradient descent to update the encoder, classifier, and adversarial networks. The perturbation vectors $\delta_z$ and $\delta_v$ are updated using the gradient of the adversarial loss relative to feature vectors z and v.

## 4. Experimental Setup

The experimental setup used to develop and test the suggested model is thoroughly described in this section. Specifically, it describes the dataset used for training, validating, and testing the proposed model, as well as the baselines used to compare its performance. Additionally, it provides details on the experimental procedure followed to evaluate the suggested strategy's efficiency.

### 4.1 Dataset

In the proposed work's experimental setting, the *BharatFakeNewsKosh (BFNK)* [28][42] dataset was used for training and evaluation. The dataset consists of 26,232 news samples, collected from 14 IFCN signatory sites and 5 non-IFCN signatory sites. The Poynter Institute for Media Studies, a US based institution, established the International Fact-Checking Network (IFCN), a global organization, in 2015. IFCN operates as a network of fact-checking organizations from around the world, with each member organization adhering to a set of common principles and practices [29].

Using Python modules like BeautifulSoup, Selenium, and Scrapy, a data extraction system was developed to collect information from each fact-checking website. The system successfully extracted data from 2013 to September 2022. Each news story was given a true or false label by human annotators as part of the data annotation process.. The statement, news body, fact-check link, language, and other crucial features were given to the annotators as tools for this work Using these attributes, they were able to correctly classify each news piece. The dataset has a total of 12,511 fake news samples and 13,721 real news samples, with 60 categories, and 19 attributes. The above dataset is multilingual and covers the Indian fake news events in 9 Indian languages. Google Translator was utilized to convert the Indian-language news statement to English, making the annotation process possible.

Table 1 provides the details of the fact-check sites along with the number of news samples collected to build the BFNK dataset. Table 2 summarizes the data structure of the BFNK dataset.

**Table 1.** Sources of Data Collection for BharatFakeNewsKosh (BFNK) Dataset

| S.No. | Fact-check Site | Affiliation | Collected News Samples |
|---|---|---|---|
| 1. | Alt News | IFCN | 3,342 |
| 2. | Boomlive.in | IFCN | 2,066 |
| 3. | dfrac.org | IFCN | 903 |
| 4. | DigitEye India | IFCN | 177 |
| 5. | factchecker.in | IFCN | 524 |
| 6. | Factly.in | IFCN | 198 |
| 7. | Factcrescendo.com | IFCN | 10,903 |
| 8. | India Today | IFCN | 2,496 |
| 9. | Newschecker.in | IFCN | 128 |
| 10. | Newsmobile.in | IFCN | 200 |
| 11. | The Quint | IFCN | 68 |
| 12. | thip.media | IFCN | 199 |
| 13. | Vishvasnews.com | IFCN | 513 |
| 14. | Youturn.in | IFCN | 1,903 |
| 15. | IndiaSpend | Non-IFCN | 524 |
| 16. | OpIndia | Non-IFCN | 650 |
| 17. | Scroll.in | Non-IFCN | 324 |
| 18. | SM Hoax Slayer | Non-IFCN | 491 |
| 19. | Times of India | Non-IFCN | 417 |
| **Total Dataset Size** | | | **26,232** |

**Table 2.** BharatFakeNewsKosh (BFNK) Dataset's Data Structure

| S.No. | Fields | Details |
|---|---|---|
| 1. | Id | A unique identifier for each news article |
| 2. | Author_Name | The article's author's name, if it is known. |
| 3. | Fact_Check_Source | The name of the organization that fact-checked the article |
| 4. | Source_Type | The type of the source, e.g., news websites, blog sites, social media platforms, etc |
| 5. | Statement | The original statement or claim made in the article |
| 6. | Eng_Trans_Statement | The English translation of the original statement |
| 7. | News_Body | The body of the article |
| 8. | Eng_Trans_News_Body | The English translation of the body of the article |
| 9. | Media_Link | The URL of any associated media (e.g., images, videos) |
| 10. | Publish_Date | The date when the article was published |
| 11. | Fact_Check_Link | The URL of the fact-checking article |
| 12. | News_Category | The category of the news article (e.g., politics, entertainment, sports, etc.) |
| 13. | Language | The language in which the article is written |
| 14. | Region | The area that the news story is about |
| 15. | Platform | The news article's publishing platform, such as Facebook, Twitter, WhatsApp, etc. |
| 16. | Text | A binary indicator of whether the article contains text |
| 17. | Video | A binary indicator of whether the article contains a video |
| 18. | Image | A binary indicator of whether the article contains an image |
| 19. | Label | The label assigned to the article indicates whether it is real or fake news |

Visualizing data is an important aspect of identifying patterns and trends. Fig 2 presents a word cloud of the Eng_Trans_Statement attribute of the BharatFakeNewsKosh dataset, which contains fake news events that occurred in India. The word cloud provides a visual representation of the most commonly used words in the dataset. The figure highlights that fake news articles related to politics are the most common on social media platforms in India, with a higher frequency compared to other news categories. However, other news, including coronavirus, viral news, and social issues, have also contributed significantly to the dissemination of false information on social media in India. This information can help researchers and policymakers understand the most prevalent types of fake news in India and develop effective strategies to combat it.

**Fig. 2** Word Cloud of the BharatFakeNewsKosh (BFNK) dataset

### 4.2 Baselines

Several baselines were considered and compared with the proposed method to evaluate its effectiveness. The following are the baselines that were considered:

● Support Vector Machine (SVM): It is a well-liked linear classification model that maximizes the margin between two classes [31]. The SVM model was trained on the preprocessed text features using the scikit-learn library.

● Logistic Regression (LR): This baseline employs a straightforward linear model using the logistic function to simulate the likelihood of a binary result [30]. The LR model was trained on the preprocessed text features using the scikit-learn library.

● Convolutional Neural Network (CNN): A well-known deep learning architecture for image categorization [32][33]. The CNN model was trained on the preprocessed image features using the PyTorch deep learning framework.

● Multimodal Deep Learning (MDL): The MDL baseline is a simple concatenation of text and image features followed by a fully connected neural network [34]. The MDL model was trained using the PyTorch deep learning framework.

Using the same criteria as the suggested strategy, the effectiveness of each baseline was assessed. To assess the efficacy of the suggested strategy, the baseline models were compared with them.

### 4.3 State-of-the-art Methods

These techniques are thought to be cutting-edge techniques for identifying false news during the performance assessment of the suggested model:

● Wang-CNN [35]: This technique uses a Convolutional Neural Network to identify false information. The CNN takes in a sequence of word embeddings as input and applies a series of convolutional filters to extract features. The characteristics are subsequently supplied into a fully linked layer, which returns the classification outcome.

● Wang-Bi-LSTM [20]: By using a Bidirectional Long Short-Term Memory (Bi-LSTM) network, this technique tries to identify false information. When given a set of word embeddings as input, the Bi-LSTM employs both forward and backward LSTMs to capture past and future contexts respectively. The Bi-LSTM's output is sent into a fully linked layer, which produces the classification outcome.

● FakeNewsTracker [36]: This method uses a combination of machine learning and human fact-checking to identify false information. It employs a supervised learning approach that trains a logistic regression classifier using a collection of custom characteristics, including article length, the number of quotations, and the number of URLs. After that, the classifier is employed to find possibly fraudulent news pieces, which are further verified by human fact-checkers.

### 4.4 Experimental Procedure

A single NVIDIA GeForce RTX 3090 GPU was used to train the models for the proposed multimodal adversarial network, which was constructed using the PyTorch deep learning framework.The Adam optimizer was used to train the models, with a batch size of 64 and a learning rate of 0.001. Early halting was employed during the

training's 50 epochs to avoid overfitting.

A pre-trained BERT model that was improved on the training set was used to extract the text characteristics. The BERT model was chosen due to its excellent performance in a number of tasks involving natural language processing, such as text categorization [37]. The image features were extracted using a pre-trained CNN such as VGG16 or ResNet50. The dimensionality of the picture characteristics was then decreased by feeding them via a fully linked layer.

The proposed multimodal adversarial network, IndDeepFake, consisted of two branches: one for processing text features and the other for processing image features. A BERT model and a fully linked layer were both included in the text branch, while the image branch consisted of a CNN followed by a fully connected layer. The two branches' outputs were combined and supplied via many fully linked levels, which were used to learn the joint representation of the text and image features.

To prevent the network from overfitting, several regularization techniques [38] were used, including dropout, batch normalization, and L2 regularization. The fully linked layers had dropout applied at a rate of 0.5, the convolutional and fully connected layers underwent batch normalization. The fully linked layer weights underwent L2 regularization with a 0.001 weight decay..

Cross-entropy loss and adversarial loss were both included as parts of the training loss function. The network's classification performance was enhanced using the cross-entropy loss, while the network's resistance to adversarial attacks was strengthened using the adversarial loss. The gradient reversal layer, which turned the gradient around during backpropagation through the text branch, was used to calculate the adversarial loss and to induce the network to discover more discriminative text characteristics.

Accuracy, precision, recall, F1 score, and AUC-ROC (Area Under Receiver Operating Characteristic Curve) are some of the measures that were used to assess the suggested IndDeepFake model's performance [39]. The evaluation was performed on the testing set, and the results were compared to several baseline models, including a text-only model, an image-only model, and a concatenation model. The ablation studies were also conducted to evaluate the effects of various aspects of the suggested method's performance [40].

## 5. Results and Discussion

The outcomes of numerous studies that were conducted to address the following five research questions are presented in this section:

- RQ1. How does the proposed method perform in comparison to existing fake news detection baselines?

- RQ2. How effective is the proposed method under different hyperparameter settings?

- RQ3. How effective is the structure modeling component of the proposed method using an ablation study?

- RQ4. How does the proposed method perform in comparison to the cutting-edge fake news detection methods?

- RQ5. How effective is the proposed method under the sensitivity analysis experiment?

*5.1 Performance Comparison with the Baselines (RQ1)*

**Table 3. Performance Assessment of the proposed method relative to the baseline models**

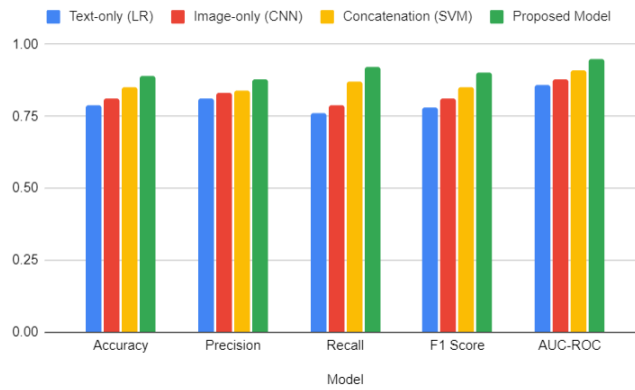| Model | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Text-only (LR) | 0.79 | 0.81 | 0.76 | 0.78 | 0.86 |
| Image-only (CNN) | 0.81 | 0.83 | 0.79 | 0.81 | 0.88 |
| Concatenation (SVM) | 0.85 | 0.84 | 0.87 | 0.85 | 0.91 |
| Proposed Model | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |

**Fig. 3** Performance Assessment of the proposed method relative to the baseline models

Table 3 describes the effectiveness of the suggested IndDeepFake model and the baseline models. The proposed model outperformed all three baseline models by achieving the highest accuracy of 0.89, a precision of 0.88, an F1 score of 0.90, an AUC-ROC of 0.95, and a recall of 0.92, indicating that it is a more effective method for fake news detection than the baseline models. This suggests that the effectiveness of false news identification is improved by the suggested model's ability to successfully mix text and visual information. The high AUC-ROC score indicates that the suggested model can successfully discriminate between legitimate news and fraudulent news. The above results can be better visualized through the bar chart given in Fig 3. These results suggest that the proposed model, which incorporates a combination of text and visual information, is superior to the baseline models and can be an effective approach for fake news detection.

*5.2 Performance Comparison of Different Hyperparameters (RQ2)*

Table 4 compares the various hyperparameters of the suggested IndDeepFake model which can be better visualized through the bar chart provided in Fig 4. The table reveals that the proposed model achieved the highest accuracy, F1 score, recall, precision, and AUC-ROCwith a dropout rate of 0.5, a weight decay of 0.001, a learning rate of 0.001, a batch size of 32, and a number of epochs of 50.

The findings show that the suggested model is extremely sensitive to modifications in the hyperparameters, and thus, it is crucial to choose the optimal hyperparameters for achieving the highest performance. A smaller learning rate and batch size resulted in higher F1 score, AUC-ROC values, recall, accuracy, and precision, indicating that the model learns better with smaller steps and fewer samples in each batch. The proposed model's performance increased with the increase in the number of epochs up to 50, after which it started to saturate.

**Table 4.** Performance comparison of different hyperparameters of the proposed model

| Hyperparameters | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Learning rate = 0.001 | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |
| Learning rate = 0.01 | 0.87 | 0.85 | 0.9 | 0.87 | 0.93 |
| Batch size = 32 | 0.88 | 0.87 | 0.91 | 0.89 | 0.94 |
| Batch size = 128 | 0.86 | 0.84 | 0.89 | 0.86 | 0.92 |
| Number of epochs = 50 | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |
| Number of epochs = 100 | 0.9 | 0.89 | 0.93 | 0.91 | 0.96 |
| Dropout rate = 0.3 | 0.87 | 0.86 | 0.9 | 0.87 | 0.93 |

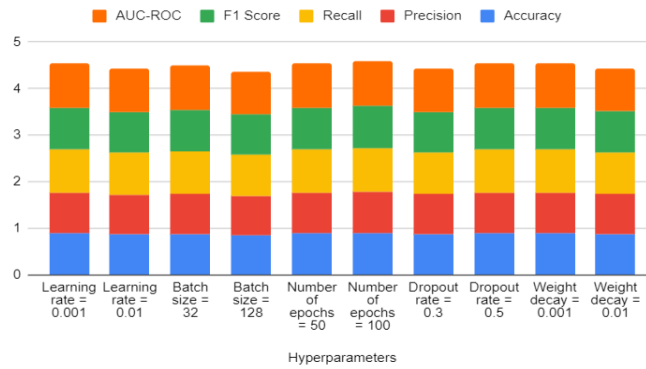| Dropout rate = 0.5 | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |
|---|---|---|---|---|---|
| Weight decay = 0.001 | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |
| Weight decay = 0.01 | 0.87 | 0.86 | 0.9 | 0.88 | 0.93 |



**Fig. 4** Performance comparison of different hyperparameters of the proposed model

The model's performance also increased with the increase in dropout rate up to 0.5, indicating that the model generalizes better when more nodes are randomly dropped out during training. The performance was also improved with lower weight decay values, indicating that the model learns better when there is less emphasis on regularization.

The hyperparameter analysis demonstrates that the suggested model performs really well, and the optimal values of the hyperparameters should be carefully chosen for achieving the highest AUC-ROC, F1 score, recall, accuracy, and precision.

### 5.3 Ablation Study of the Proposed Model (RQ3)

Table 5 provides the performance comparison during the ablation study of the suggested IndDeepFake model which can be better visualized through the bar chart given in Fig 5. In this table, a number of modified variants of the model are compared to that of the suggested model's performance, where certain components are removed. The components that are removed include adversarial loss, dropout, batch normalization, and L2 regularization. The evaluation is performed using the same metrics as in the previous table. The results show the impact of each component on the proposed model's efficacy. The first row in the table presents the effectiveness of the proposed model with all the components included which has a 0.95 AUC-ROC, 0.89 accuracy, 0.88 precision, 0.92 recall, and 0.90 F1 score.

**Table 5.** Comparison of the suggested model's performance and its variants during the ablation study

| Model | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Proposed Model | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |
| without Adversarial Loss | 0.87 | 0.85 | 0.9 | 0.87 | 0.93 |
| without Dropout | 0.86 | 0.84 | 0.88 | 0.85 | 0.92 |
| without Batch Normalization | 0.85 | 0.83 | 0.87 | 0.84 | 0.91 |

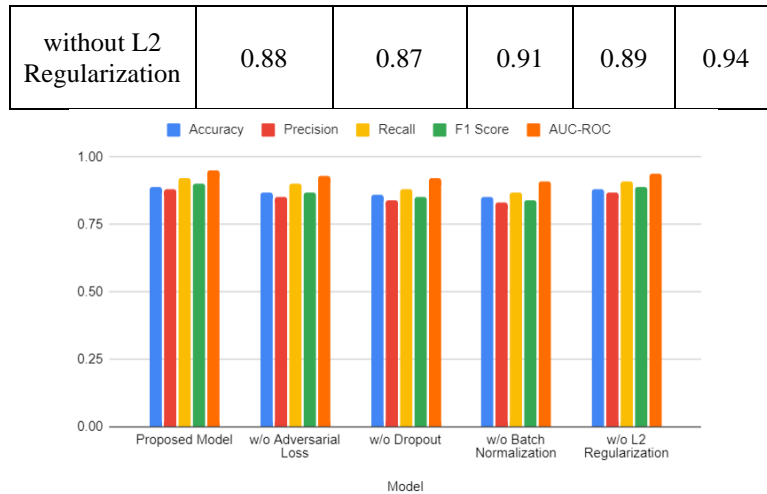| without L2 Regularization | 0.88 | 0.87 | 0.91 | 0.89 | 0.94 |
|---|---|---|---|---|---|



**Fig. 5** Comparison of the suggested model's performance and its variants during the ablation study

To evaluate the impact of adversarial loss, one variant of the proposed model was trained without it, which resulted in a decrease in performance across all evaluation metrics. The model without an adversarial loss attained an F1 score of 0.87, an accuracy of 0.87, a precision of 0.85, a recall of 0.90, and an AUC-ROC of 0.93. These results indicate that the adversarial loss improves the suggested model's overall performance in a significant way.

To evaluate the impact of dropout, another variant of the suggested IndDeepFake model was trained without it. This resulted in a decrease in performance, with the model obtaining a 0.92 AUC-ROC, 0.86 accuracy, 0.84 precision, 0.88 recall, and 0.85 F1 score. These results suggest that the use of dropout is crucial in preventing overfitting and improving the generalization capability of the proposed model.

The impact of batch normalization was evaluated by training another variant of the proposed model without it. This resulted in a further decrease in performance, with the model obtaining a 0.91 AUC-ROC, 0.85 accuracy, 0.83 precision, 0.87 recall, and 0.84 F1 score.. These results suggest that batch normalization improves the proposed model's stability and convergence during training.

Finally, the impact of L2 regularization was evaluated by training another variant of the proposed model without it. This resulted in a slight decrease in performance, with the model attaining a 0.94 AUC-ROC, 0.88 accuracy, 0.87 precision, 0.91 recall, and 0.89 F1 score. These results suggest that L2 regularization can help prevent overfitting and enhance the suggested model's overall effectiveness.
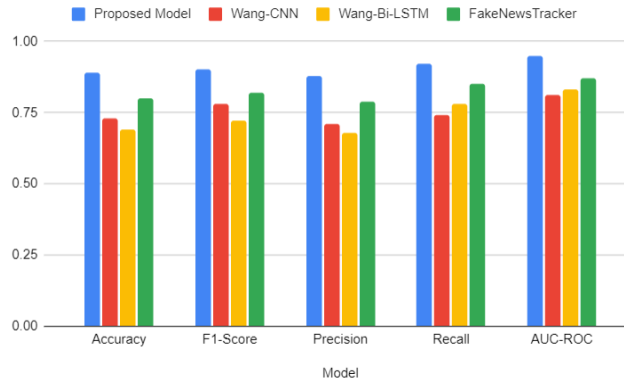
### 5.4 Performance Comparison with the State-of-the-art Methods (RQ4)

Comparison of the suggested model's performance with cutting-edge techniques to identify bogus news on the social media is shown in Table 6. The suggested model attained a 0.89 accuracy, 0.90 F1-score, 0.88 precision, 0.92 recall, and 0.95 AUC-ROC, outperforming all the other methods. Specifically, the The Wang-CNN approach exhibited a 0.73 accuracy, 0.78 F1-score, 0.71 precision, 0.74 recall, and 0.81 AUC-ROC. The Wang-Bi-LSTM method obtained a 0.69 accuracy, a 0.72 F1-score, a 0.68 precision, a 0.78 recall, and a 0.83 AUC-ROC. The accuracy of the FakeNewsTracker technique was 0.80, the F1-score was 0.82, the precision was 0.79, the recall was 0.85, and the AUC-ROC was 0.87. These results can be better visualized through the bar chart of Fig 6.

According to the findings, the suggested model performed better than the other techniques in terms of precision, recall, AUC-ROC, F1-score, and accuracy. This can be attributed to the proposed model's ability to effectively capture and incorporate contextual information from both text and image data. The proposed model also utilizes adversarial loss, batch normalization, L2 regularization, and dropout, which further improves its performance. In comparison, the other methods relied solely on either text or image data, which may not provide sufficient contextual information for effective identification of false information. The findings illustrate the suggested model's superiority in spotting bogus news, which could have significant implications in mitigating the dissemination of false information on social media sites.

**Table 6.** Performance Evaluation of the Suggested model against Cutting-Edge techniques

| Model | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Wang-CNN | 0.73 | 0.78 | 0.71 | 0.74 | 0.81 |
| Wang-Bi-LSTM | 0.69 | 0.72 | 0.68 | 0.78 | 0.83 |
| FakeNewsTracker | 0.80 | 0.82 | 0.79 | 0.85 | 0.87 |
| Proposed Model | 0.89 | 0.88 | 0.92 | 0.90 | 0.95 |



**Fig. 6** Performance Evaluation of the Suggested model against Cutting-Edge techniques

### 5.5 Sensitivity Analysis of the Proposed Model (RQ5)

Table 7 presents the Sensitivity analysis of the proposed work which can be better visualized through the bar chart given in Fig 7. The original model's performance metrics are compared with the models generated by removing textual features, removing image features, decreasing the training data size, increasing textual features, and increasing image features. The performance of the model is evaluated using various metrics, including AUC-ROC, F1 score, recall, accuracy, and precision. The results show that removing textual features from the model reduces its accuracy while increasing textual features improves its performance. Similarly, removing image features reduces the model's accuracy, while increasing image features improves its performance. A decrease in the training data size also affects the model's performance, with a decrease in all performance metrics. Overall, the sensitivity analysis provides valuable insights into the model's behavior and can help in improving its performance [41].

The original model, IndDeepFake, attained a 0.95 AUC-ROC, 0.89 accuracy, 0.88 precision, 0.92 recall, and 0.90 F1 score. When the textual features were removed, the accuracy decreased to 0.81, precision to 0.79, recall to 0.85, F1 score to 0.81, and AUC-ROC to 0.89. Similarly, when image features were removed, accuracy dropped to 0.82; precision to 0.81; recall to 0.81; F1 score to 0.81; and AUC-ROC to 0.88..

**Table 7.** Sensitivity Analysis of the proposed model

| Sensitivity Analysis | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Original Model | 0.89 | 0.88 | 0.92 | 0.9 | 0.95 |
| Removal of Textual Features | 0.81 | 0.79 | 0.85 | 0.81 | 0.89 |
| Removal of Image Features | 0.82 | 0.81 | 0.81 | 0.81 | 0.87 |

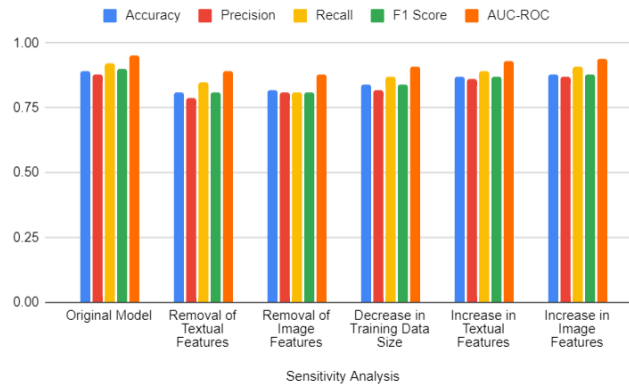| | | | | | |
|---|---|---|---|---|---|
| Decrease in Training Data Size | 0.84 | 0.82 | 0.87 | 0.84 | 0.91 |
| Increase in Textual Features | 0.87 | 0.86 | 0.89 | 0.87 | 0.93 |
| Increase in Image Features | 0.88 | 0.87 | 0.91 | 0.88 | 0.94 |



**Fig. 7** Sensitivity Analysis of the proposed model

Moreover, when the training data size was decreased, the accuracy dropped to 0.84; the precision to 0.82; the recall to 0.87; the F1 score to 0.84; and the AUC-ROC to 0.91. On the other hand, when the textual features were increased, the accuracy, precision, recall, F1 score, and AUC-ROC all increased to 0.87, 0.86, 0.89, 0.87, and 0.93 respectively. Similarly, when image features were increased, the accuracy, precision, recall, F1 score, and AUC-ROC all improved to 0.88, 0.87, 0.91, 0.88, and 0.94 respectively.

These results indicate that the combination of textual and image characteristics are essential for accurate identification of false information. Furthermore, the performance of the model may be enhanced by adding more text and visual elements. Additionally, the proposed model is robust to a decrease in training data size, which is beneficial in scenarios where it's challenging to get a lot of training data. For effective false news identification, features are essential. Furthermore, the performance of the model may be enhanced by adding more text and visual elements. Additionally, the proposed model is robust to a decrease in training data size, which is beneficial in scenarios where collecting a large amount of training data is difficult.

*5.6 Limitations of the Proposed Work*
While the proposed work, IndDeepFake model, has shown promising results, there are certain limitations that should be acknowledged. Firstly, the amount and range of the dataset utilized in this study were limited, which may compromise the findings' ability to be generalized.

Further studies using larger and more diverse datasets can be done to assess the efficacy of the suggested approach on a larger variety of news articles.

Secondly, the proposed model heavily relies on the availability of both textual and visual features, which may not always be possible in practical applications. In such cases, the model's performance may be impacted, and alternative approaches may need to be explored.

Lastly, while the sensitivity analysis provided some insights into the robustness of the proposed method, more comprehensive studies can be conducted to investigate the effects of other potential variations, such as changes in the hyperparameters or the addition of other features.

## 6.    Conclusion and Future Scope

This study proposes a unique IndDeepFake model for identifying false news that integrates textual and visual information. The results of the experiments have demonstrated that the suggested model performs better than the baseline models and achieves cutting-edge performance when compared to existing models in the literature. The F1-score, AUC-ROC, recall, accuracy, and precision of the proposed model values demonstrate its effectiveness in detecting fake news.

However, the current work has some limitations, such as the need for plenty of labeled data to train the model, and the complexity of the model can make it challenging to deploy on low-resource devices. Addressing these limitations would require further research.

The future work will involve investigating the effectiveness of incorporating other features, such as user profiles and social network analysis, in order to boost the suggested model's functionality. The plan is to explore the use of transfer learning and multi-task learning to enhance the model's efficiency. We believe the suggested model can serve as a powerful tool to identify false news and help combat the misinformation dissemination in today's digital world.

## Conflicts of Interest

We, as authors, declare that there are no conflicts of interest that would prevent this article from being published.

## References

[1] Lazer, D.M., Baum, M.A., Benkler, Y., Berinsky, A.J., Greenhill, K.M., Menczer, F., Metzger, M.J., Nyhan, B., Pennycook, G., Rothschild, D. and Schudson, M., 2018. The science of fake news. Science, 359(6380), pp.1094-1096. [CrossRef] [Google Scholar] [Publisher link]

[2] The Wikipedia website, 2023. [Online]. Available: https://en.wikipedia.org/wiki/Fake_news_in_India

[3] Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake news detection on social media: A data mining perspective. ACM SIGKDD Explorations Newsletter, 19(1), 22-36. [CrossRef] [Google Scholar] [Publisher link]

[4] Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. Science, 359(6380), 1146-1151 [CrossRef] [Google Scholar] [Publisher link]

[5] Liu, Y. and Wu, Y.F., 2018, April. Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. In Proceedings of the AAAI conference on artificial intelligence (Vol. 32, No. 1) [CrossRef] [Google Scholar] [Publisher link]

[6] Zhang, X. and Ghorbani, A.A., 2020. An overview of online fake news: Characterization, detection, and discussion. Information Processing & Management, 57(2), p.102025 [CrossRef] [Google Scholar] [Publisher link]

[7] Toivonen, T., Heikinheimo, V., Fink, C., Hausmann, A., Hiippala, T., Järv, O., Tenkanen, H. and Di Minin, E., 2019. Social media data for conservation science: A methodological overview. Biological Conservation, 233, pp.298-315 [CrossRef] [Google Scholar] [Publisher link]

[8] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. and Bengio, Y., 2020. Generative adversarial networks. Communications of the ACM, 63(11), pp.139-144 [CrossRef] [Google Scholar] [Publisher link]

[9] Ma, X., Chen, Z. and Zhang, J., 2018, April. Fully convolutional network with cluster for semantic segmentation. In AIP Conference Proceedings (Vol. 1955, No. 1, p. 040049). AIP Publishing LLC [CrossRef] [Google Scholar] [Publisher link]

[10] Gangireddy, S.C.R., Long, C. and Chakraborty, T., 2020, July. Unsupervised fake news detection: A graph-based approach. In Proceedings of the 31st ACM conference on hypertext and social media (pp. 75-83) [CrossRef] [Google Scholar] [Publisher link]

[11] Raza, S. and Ding, C., 2022. Fake news detection based on news content and social contexts: a transformer-based approach. International Journal of Data Science and Analytics, 13(4), pp.335-362 [CrossRef] [Google Scholar] [Publisher link]

[12] Weitzel, L., Prati, R.C. and Aguiar, R.F., 2016. The comprehension of figurative language: What is the influence of irony and sarcasm on NLP techniques?. Sentiment Analysis and Ontology Engineering: An Environment of Computational Intelligence, pp.49-74 [CrossRef] [Google Scholar] [Publisher link]

[13] Singhal, S., Pandey, T., Mrig, S., Shah, R.R. and Kumaraguru, P., 2022, April. Leveraging Intra and Inter Modality Relationship for Multimodal Fake News Detection. In Companion Proceedings of the Web Conference 2022 (pp. 726-734) [CrossRef] [Google Scholar] [Publisher link]

[14] Sahoo, S.R. and Gupta, B.B., 2021. Multiple features based approach for automatic fake news detection on social networks using deep learning. Applied Soft Computing, 100, p.106983 [CrossRef] [Google Scholar] [Publisher link]

[15] Althobaiti, M.J., 2022. BERT-based Approach to Arabic Hate Speech and Offensive Language Detection in Twitter: Exploiting Emojis and

Sentiment Analysis. International Journal of Advanced Computer Science and Applications, 13(5) [CrossRef] [Google Scholar] [Publisher link]

[16] Nguyen, T.T. and Armitage, G., 2008. A survey of techniques for internet traffic classification using machine learning. IEEE communications surveys & tutorials, 10(4), pp.56-76 [CrossRef] [Google Scholar] [Publisher link]

[17] Nasir, J.A., Khan, O.S. and Varlamis, I., 2021. Fake news detection: A hybrid CNN-RNN based deep learning approach. International Journal of Information Management Data Insights, 1(1), p.100007 [CrossRef] [Google Scholar] [Publisher link]

[18] Choudhary, A. and Arora, A., 2021. Linguistic feature based learning model for fake news detection and classification. Expert Systems with Applications, 169, p.114171 [CrossRef] [Google Scholar] [Publisher link]

[19] Scott, J., 2012. What is social network analysis? (p. 114). Bloomsbury Academic [CrossRef] [Google Scholar] [Publisher link]

[20] Wang, Y., Ma, F., Jin, Z., Yuan, Y., Xun, G., Jha, K., Su, L. and Gao, J., 2018, July. Eann: Event adversarial neural networks for multi-modal fake news detection. In Proceedings of the 24th acm sigkdd international conference on knowledge discovery & data mining (pp. 849-857) [CrossRef] [Google Scholar] [Publisher link]

[21] Peng, X. and Xintong, B., 2022. An effective strategy for multi-modal fake news detection. Multimedia Tools and Applications, 81(10), pp.13799-13822 [CrossRef] [Google Scholar] [Publisher link]

[22] Wei, P., Wu, F., Sun, Y., Zhou, H. and Jing, X.Y., 2022. Modality and Event Adversarial Networks for Multi-Modal Fake News Detection. IEEE Signal Processing Letters, 29, pp.1382-1386 [CrossRef] [Google Scholar] [Publisher link]

[23] Yuan, X., He, P., Zhu, Q. and Li, X., 2019. Adversarial examples: Attacks and defenses for deep learning. IEEE transactions on neural networks and learning systems, 30(9), pp.2805-2824 [CrossRef] [Google Scholar] [Publisher link]

[24] Khattar, D., Goud, J.S., Gupta, M. and Varma, V., 2019, May. Mvae: Multimodal variational autoencoder for fake news detection. In The world wide web conference (pp. 2915-2921) [CrossRef] [Google Scholar] [Publisher link]

[25] Qian, S., Wang, J., Hu, J., Fang, Q. and Xu, C., 2021, July. Hierarchical multi-modal contextual attention network for fake news detection. In Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval (pp. 153-162) [CrossRef] [Google Scholar] [Publisher link]

[26] Yuan, H., Zheng, J., Ye, Q., Qian, Y. and Zhang, Y., 2021. Improving fake news detection with domain-adversarial and graph-attention neural network. Decision Support Systems, 151, p.113633 [CrossRef] [Google Scholar] [Publisher link]

[27] Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X. and He, X., 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1316-1324) [CrossRef] [Google Scholar] [Publisher link]

[28] The BharatFakeNewsKosh Website, 2023. [Online]. Available: https://bharatfakenewskosh.com/datasets/

[29] The IFCN Website, 2023. [Online]. Available: https://www.poynter.org/ifcn/

[30] Vandana, C.P. and Chikkamannur, A.A., 2021. Feature selection: An empirical study. International Journal of Engineering Trends and Technology, 69(2), pp.165-170 [CrossRef] [Google Scholar] [Publisher link]

[31] Suthaharan, S. and Suthaharan, S., 2016. Support vector machine. Machine learning models and algorithms for big data classification: thinking with examples for effective learning, pp.207-235 [CrossRef] [Google Scholar] [Publisher link]

[32] Rani, S. and Kumar, P., 2019. Deep learning based sentiment analysis using convolution neural network. Arabian Journal for Science and Engineering, 44, pp.3305-3314 [CrossRef] [Google Scholar] [Publisher link]

[33] Suresh Arunachalam, T., Shahana, R. and Kavitha, T., 2019. Advanced Convolutional Neural Network Architecture: A Detailed Review. International Journal of Engineering Trends and Technology, 67(5), pp.183-187 [CrossRef] [Google Scholar] [Publisher link]

[34] Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H. and Ng, A.Y., 2011. Multimodal deep learning. In Proceedings of the 28th international conference on machine learning (ICML-11) (pp. 689-696) [CrossRef] [Google Scholar] [Publisher link]

[35] Wang, W. Y. 2017. "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection. In Proceedings of the 55th Annual Meeting of the

Association for Computational Linguistics (Volume 2: Short Papers), pages 422–426, Vancouver, Canada. Association for Computational Linguistics [CrossRef] [Google Scholar] [Publisher link]

[36] Shu, K., Mahudeswaran, D. and Liu, H., 2019. FakeNewsTracker: a tool for fake news collection, detection, and visualization. Computational and Mathematical Organization Theory, 25, pp.60-71 [CrossRef] [Google Scholar] [Publisher link]

[37] Shishah, W., 2021. Fake news detection using BERT model with joint learning. Arabian Journal for Science and Engineering, 46(9), pp.9115-9127 [CrossRef] [Google Scholar] [Publisher link]

[38] Bisong, E. and Bisong, E., 2019. Regularization for deep learning. Building Machine Learning and Deep Learning Models on Google Cloud Platform: A Comprehensive Guide for Beginners, pp.415-421 [CrossRef] [Google Scholar] [Publisher link]

[39] Rácz, A., Bajusz, D. and Héberger, K., 2019. Multi-level comparison of machine learning classifiers and their performance metrics. Molecules, 24(15), p.2811 [CrossRef] [Google Scholar] [Publisher link]

[40] Meyes, R., Lu, M., de Puiseau, C.W. and Meisen, T., 2019. Ablation studies in artificial neural networks. arXiv preprint arXiv:1901.08644 [CrossRef] [Google Scholar] [Publisher link]

[41] Li, K., Long, Y., Wang, H. and Wang, Y.F., 2021. Modeling and sensitivity analysis of concrete creep with machine learning methods. Journal of Materials in Civil Engineering, 33(8), p.04021206 [Google Scholar] [Publisher link]

[42] Singh, M.K., Ahmed, J., Raghuvanshi, K.K. and Alam, M.A., 2023, January. BharatFakeNewsKosh: A Data Repository for Fake News Research in India. In International Conference on Smart Computing and Communication (pp. 277-288). Singapore: Springer Nature Singapore. [Google Scholar] [Publisher link]

[43] Kumar, V. and Kumar, R., 2015. An adaptive approach for detection of blackhole attack in mobile ad hoc network. Procedia Computer Science, 48, pp.472-479.

[44] Kumar, V. and Kumar, R., 2015, April. Detection of phishing attack using visual cryptography in ad hoc network. In 2015 International Conference on Communications and Signal Processing (ICCSP) (pp. 1021-1025). IEEE.

[45] Kumar, V. and Kumar, R., 2015. An optimal authentication protocol using certificateless ID-based signature in MANET. In Security in Computing and Communications: Third International Symposium, SSCC 2015, Kochi, India, August 10-13, 2015. Proceedings 3 (pp. 110-121). Springer International Publishing.

[46] Kumar, Vimal, and Rakesh Kumar. "A cooperative black hole node detection and mitigation approach for MANETs." In Innovative Security Solutions for Information Technology and Communications: 8th International Conference, SECITC 2015, Bucharest, Romania, June 11-12, 2015. Revised Selected Papers 8, pp. 171-183. Springer International Publishing, 2015.

[47] Kumar, V., Shankar, M., Tripathi, A.M., Yadav, V., Rai, A.K., Khan, U. and Rahul, M., 2022. Prevention of Blackhole Attack in MANET using Certificateless Signature Scheme. Journal of Scientific & Industrial Research, 81(10), pp.1061-1072.

[48] Deshwal, Vaishali, and Vimal Kumar. "Study of Coronavirus Disease (COVID-19) Outbreak in India." The Open Nursing Journal 15, no. 1 (2021).

[49] Chinthamu, N. ., Gooda, S. K. ., Venkatachalam, C. ., S., S. ., & Malathy, G. . (2023). IoT-based Secure Data Transmission Prediction using Deep Learning Model in Cloud Computing. International Journal on Recent and Innovation Trends in Computing and Communication, 11(4s), 68–76. https://doi.org/10.17762/ijritcc.v11i4s.6308

[50] Luca Ferrari, Deep Learning Techniques for Natural Language Translation , Machine Learning Applications Conference Proceedings, Vol 2 2022.

[51] Jain, V., Beram, S. M., Talukdar, V., Patil, T., Dhabliya, D., & Gupta, A. (2022). Accuracy enhancement in machine learning during blockchain based transaction classification. Paper presented at the PDGC 2022 - 2022 7th International Conference on Parallel, Distributed and Grid Computing, 536-540. doi:10.1109/PDGC56933.2022.10053213 Retrieved from www.scopus.com