

AN INTERPRETABLE DEEP LEARNING FRAMEWORK FOR DETECTING MISINFORMATION IN DIGITAL PLATFORMS

[1] Mazid Mahmood, M.Tech Research Scholar, Department of Computer Science and Engineering, Integral University

[2] Dr. Nudrat Fatima, Supervisor, Department of Computer Science and Engineering, Integral University, Lucknow

[3] Dr. Sifatullah Siddiqi, Co-Supervisor, Department of Computer Science and Engineering, Integral University, Lucknow

ABSTRACT: The proliferation of misinformation on digital platforms poses significant challenges to society, affecting public opinion, health, politics, and security. While deep learning models have demonstrated high accuracy in detecting fake news and misinformation, their black-box nature limits transparency and trustworthiness. This study proposes an interpretable deep learning framework that not only identifies misinformation with high precision but also provides human-understandable explanations for its predictions. By integrating attention mechanisms, explainable AI (XAI) techniques, and multi-modal data analysis, the proposed framework enhances both detection performance and interpretability. Experimental evaluation on benchmark datasets showcases the model's superiority in accuracy, explainability, and generalizability. In the digital age, the rapid spread of misinformation challenges the credibility of news consumption, prompting this study to propose an Integrated Hybrid Deep Learning AI (IHDLAI) framework that leverages advanced deep learning and AI-driven natural language processing techniques for effective misinformation detection and mitigation; the framework consists of three phases: Data Collection and Pre-processing, where datasets from verified news sources,

social media, and fact-checking websites are cleaned, tokenized, normalized, and feature-extracted to build a robust dataset; Model Development and Training, where deep learning models such as CNNs, RNNs with LSTM, and transformer-based models like BERT are trained, optimized, and enhanced with explainable AI methods to ensure transparency; and Evaluation and Deployment, where model performance is rigorously assessed through metrics like accuracy, precision, recall, and F1-score, with the best model achieving 93% accuracy, 92% precision, 91% recall, and a 91.5% F1-score, demonstrating the IHDLAI framework's effectiveness in combating misinformation and providing a reliable, interpretable solution to enhance the integrity of digital news platforms.

Keywords: Misinformation Detection, Deep Learning, Explainable AI (XAI), Fake News, Digital Platforms, Interpretability, Multi-modal Analysis.

1. INTRODUCTION

In today's digital era, the unprecedented proliferation of misinformation across online platforms poses significant challenges to the authenticity and reliability of news consumption. The rapid dissemination of false information not only misguides public opinion but also threatens social trust and democratic processes. Despite efforts by fact-checking agencies and social media platforms, manual verification is time-consuming and insufficient to curb the widespread impact of fake news. To address this pressing issue, the integration of Artificial Intelligence (AI) and Deep Learning (DL) has emerged as a promising solution for automated misinformation detection. This study proposes an Integrated Hybrid Deep

Learning AI (IHDLAI) framework, designed to accurately detect and mitigate the spread of misinformation by leveraging advanced Natural Language Processing (NLP) techniques, interpretable deep learning models, and hybrid architectures. The framework focuses on ensuring both high detection accuracy and model explainability, thereby enhancing transparency and user trust. Through a structured approach involving data collection, pre-processing, model development, and evaluation, the proposed IHDLAI framework aims to provide a robust and scalable solution for safeguarding the credibility of digital information ecosystems. Misinformation has become a pervasive issue on digital platforms such as social media, news websites, and messaging apps [1] [2]. The rapid dissemination of false or misleading information can influence public perception, endanger public health, and destabilize political processes. Traditional methods of misinformation detection rely on manual fact-checking, which is time-consuming and insufficient to handle the vast volume of data generated daily.

Deep learning models, particularly Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformers, have shown promising results in automating misinformation detection. However, these models often act as "black boxes," making it difficult for users and stakeholders to understand why a particular piece of content is classified as misinformation.

To address this gap, this research proposes an interpretable deep learning framework that combines high detection accuracy with model explainability. The framework leverages advanced techniques in Explainable AI (XAI), such as

attention visualization, Layer-wise Relevance Propagation (LRP), and SHAP (SHapley Additive exPlanations), to provide insights into model decisions. The inclusion of multi-modal data (text, images, metadata) further enhances detection capabilities.

The phenomenon of fake news is not a recent development; its origins can be traced back to the early days of print media where misinformation and sensationalized content, often termed "yellow journalism," were utilized to attract readership and influence public opinion. In these pre-digital times, misinformation dissemination was slow and controllable, with traditional gatekeepers such as editors and journalists serving as the primary filters of news content. However, the advent of the digital revolution and the rapid rise of social media platforms have radically altered the information landscape, resulting in an unprecedented scale and speed of misinformation dissemination [2] [3].

With the proliferation of smartphones, inexpensive internet access, and the emergence of platforms like Facebook, Twitter, WhatsApp, Instagram, and YouTube, the dynamics of news consumption and distribution have undergone a significant transformation. Unlike traditional media outlets, social media allows users to create, share, and consume content instantly without any rigorous fact-checking or editorial oversight. This decentralized and unregulated environment has created fertile ground for the rapid spread of misinformation, often referred to as fake news.

Definition and Nature of Fake News

Fake news can be defined as fabricated, misleading, or manipulated information that is presented as genuine news with the intent to deceive, influence, or misinform readers. It can take various forms, including:

Fabricated content: Completely false information designed to deceive.

Manipulated content: Genuine information distorted to mislead.

Misleading headlines: Sensational or clickbait titles that misrepresent the actual content. Satire and parody: Humorous or exaggerated content often misunderstood as factual. False context: Genuine content shared with incorrect contextual information. While satire and parody may not intend to deceive, the lines between intentional misinformation and unintentional misinformation often blur, especially when content goes viral without proper contextual understanding.

The Role of Social Media in Fake News Propagation

Social media platforms serve as catalysts for the rapid spread of misinformation due to their unique features:

Viral Sharing Mechanisms: Features like retweets, shares, and forwards enable instant dissemination to wide audiences.

Echo Chambers: Algorithms that curate content based on user preferences create information silos, reinforcing existing beliefs and reducing exposure to diverse viewpoints. Anonymity and Lack of Accountability: Users can easily create fake accounts and disseminate false information without repercussions.

Low Barriers to Entry: Anyone with internet access can publish content, blurring the lines between credible journalism and user-generated misinformation.

Emotional Content Amplification: Fake news often leverages emotional appeals, sensationalism, and fear-mongering to increase engagement and virality.

Studies have shown that fake news spreads faster than true news, as it often contains surprising, novel, or emotionally charged content that compels users to share. This creates a cyclical propagation model, where misinformation perpetuates rapidly, influencing public perception before fact-checking mechanisms can respond.

Impacts of Fake News

The societal impacts of fake news are multifaceted and profound:

Political Manipulation: Fake news has been weaponized in political campaigns to mislead voters, polarize opinions, and manipulate electoral outcomes.

Public Health Risks: During crises such as the COVID-19 pandemic, misinformation about treatments, vaccines, and preventive measures led to widespread panic, vaccine hesitancy, and public health challenges.

Social Unrest: False narratives can incite violence, communal disharmony, and mass protests based on fabricated incidents [3] [4] [5].

Economic Consequences: Fake news targeting companies or financial markets can influence stock prices, investor decisions, and economic stability.

Erosion of Trust: Continuous exposure to misinformation undermines public trust in legitimate news sources, scientific institutions, and democratic processes.

Existing Approaches for Fake News Detection

Over the past decade, researchers have explored various computational methods for detecting fake news. Broadly, these approaches can be categorized into:

1. Knowledge-Based Approaches

These involve verifying the factual accuracy of content by cross-referencing it with reliable knowledge bases or fact-checking repositories. While effective, these methods are often manual, time-consuming, and lack scalability in real-time scenarios.

2. Style-Based Approaches

These focus on identifying linguistic and stylistic patterns typical of fake news articles, such as excessive use of superlatives, emotionally charged language, or unusual writing styles. However, they may struggle with satirical content or sophisticated fake news that mimics legitimate news writing styles.

3. Propagation-Based Approaches

These analyze the diffusion patterns of information on social networks, identifying anomalies in how fake news spreads compared to legitimate news. Graph-based models and network analysis are commonly used, but these methods are computationally intensive and depend on access to complete network data.

4. Machine Learning and Deep Learning Approaches

Recent advancements in machine learning (ML) and deep learning (DL) have revolutionized fake news detection by enabling automated analysis of large-scale data. Techniques such as Naïve Bayes, Support Vector Machines (SVM), Random

Forests, and Gradient Boosting have been applied to classify news articles as real or fake based on feature sets derived from text, metadata, or user behavior.

Deep learning models, particularly Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) cells, and transformer-based models like BERT (Bidirectional Encoder Representations from Transformers), have demonstrated superior performance by capturing complex linguistic patterns, semantic relationships, and contextual nuances in textual data.

However, despite their high accuracy, many deep learning models lack interpretability, functioning as "black boxes" that provide predictions without clear explanations. This opaqueness limits their adoption in sensitive domains where transparency and accountability are paramount.

The Need for Interpretable AI in Fake News Detection

In the era of Explainable Artificial Intelligence (XAI), the emphasis is shifting towards models that not only perform well but also offer insights into their decision-making processes. Interpretable AI aims to bridge the gap between accuracy and transparency, providing:

Feature Importance Visualization: Highlighting key words, phrases, or patterns that influenced the model's prediction.

Model-Agnostic Explanations: Techniques like LIME (Local Interpretable Model-Agnostic Explanations) and SHAP (SHapley Additive exPlanations) offer post-hoc explanations for model outputs.

Attention Mechanisms: Neural attention layers focus on the most relevant parts of the input data, enhancing interpretability while maintaining performance.

Incorporating interpretability into fake news detection models is essential for building user trust, facilitating fact-checking, and enabling informed decision-making by stakeholders [6] [7] [8].

Technological Challenges

Despite significant advancements, several challenges persist in the domain of fake news detection:

Multilingual and Cross-Cultural Variations: Fake news manifests differently across languages and cultures, necessitating models that can generalize beyond English-centric datasets.

Adversarial Content Creation: Malicious actors continuously evolve their strategies to bypass detection algorithms, creating a cat-and-mouse dynamic.

Real-Time Detection Requirements: Effective mitigation of misinformation requires models that can operate at scale and in real-time, processing vast volumes of data efficiently [8] [9] [10].

2. RELATED WORKS

The proliferation of fake news, especially across news outlets and social media platforms, has spurred extensive research into detection techniques leveraging Artificial Intelligence (AI) and Machine Learning (ML). A wide spectrum of literature focuses on enhancing detection accuracy, model transparency, and

domain adaptability. Recent advancements underline hybrid models, explainable AI (XAI), and multimodal analysis as key focus areas.

Recent Advancements in Fake News Detection

Hashmi et al. (2024) introduced a hybrid framework integrating FastText embeddings with Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) models. Their approach achieved commendable accuracy and F1 scores across diverse datasets. Moreover, they explored transformer-based models like BERT and RoBERTa, embedding explainable AI techniques to demystify model decisions and enhance user trust [10] [11].

Comito et al. (2023) provided a comprehensive survey on deep learning methods for **multimodal fake news detection**. Their study emphasized feature fusion strategies, data heterogeneity, and challenges in integrating textual, visual, and social signals. The authors stressed the need for enhanced model interpretability and proposed future directions focusing on **cross-domain adaptability** and **explainability**.

Yuan et al. (2023) reviewed fake news detection technologies from a multidisciplinary lens, incorporating insights from linguistics, psychology, and communication sciences. They advocated for sustainable **human-machine interaction systems** for information dissemination, thereby highlighting the social and cognitive aspects of misinformation spread.

Joshi et al. (2023) employed Domain Adversarial Neural Networks (DANN) for **cross-platform fake news detection**. By integrating Local Interpretable

Model-Agnostic Explanations (LIME), their approach improved both detection accuracy and interpretability, particularly for COVID-19 misinformation across diverse social media environments [11] [12] [13].

Stance Detection and Evidence-Based Models

An emerging research trajectory involves leveraging **stance detection** for fake news classification. **Yuan et al. (2023)** demonstrated a novel methodology combining stance labels and causal inference using propensity score matching. Their findings revealed a negative correlation between stance consistency and fake news classification, enhancing the contextual understanding of misinformation.

Xu and Kechadi (2024) introduced a **fuzzy deep learning model** evaluated on the LIAR dataset, proposing an enhanced dataset version, LIAR2, to address data limitations. Their work achieved state-of-the-art results, focusing on improving data quality for better detection efficacy.

Ali et al. (2023) designed a **Web-Informed-Augmented Fake News Detection Model**, utilizing CNNs and deep autoencoders. By augmenting data from credible web sources, their approach significantly enhanced model robustness and accuracy in challenging scenarios.

Explainability, Social Explanations, and Graph-Based Approaches

To address model opacity, **Gong et al. (2024)** introduced the concept of **social explanations** in XAI, advocating for interdisciplinary collaboration to combat misinformation effectively. They stressed the socio-technical implications of model interpretability in real-world deployments.

Schütz (2023) combined domains such as fake news, hate speech, and propaganda, leveraging transfer learning and knowledge graphs. Their multi-domain approach, coupled with XAI methods, aimed to create interpretable and generalizable models.

Ayetiran and Özgöbek (2024) reviewed deep learning techniques for multimodal fake news detection, discussing data fusion challenges, evolving architectures, and future prospects in the integration of harmful language detection.

Siam (2024) emphasized the superiority of **advanced mixed models** over traditional classifiers, stressing the importance of architecture design tailored to social media's dynamic misinformation landscape [13] [14] [15] [16].

Zaki et al. (2024) proposed a **graph-based node embedding approach** for fake review detection, achieving high accuracy using Node2Vec embeddings. Their model, applied to datasets like Deceptive Opinion Spam Corpus and YelpChi, integrated XAI for transparent predictions.

Xia et al. (2023) developed a hybrid architecture combining CNN, BiLSTM, and Attention mechanisms, enhancing detection accuracy while providing insights into feature importance.

Multimodal, Multilingual, and Emerging Trends

Nadeem et al. (2023) proposed a Stylometric and Semantic Similarity-Oriented Framework for multimodal fake news detection, demonstrating superior performance across benchmark datasets.

Alghamdi et al. (2024) presented a multilingual detection model leveraging hybrid summarization and mBERT for classification. Their approach effectively reduced data redundancy while maintaining critical information, setting new performance standards.

Tufchi et al. (2023) offered a comprehensive survey on disinformation detection, analyzing the lifecycle of fake news and evaluating various classification models, highlighting gaps and mitigation strategies.

Vishnupriya et al. (2024) explored innovative trends like **user behavior analysis** and **collaborative cybersecurity strategies**, advocating for interdisciplinary cooperation to tackle deception in digital ecosystems [16] [17] [18] [19] [20].

Identified Research Gaps

Despite these advancements, several critical research gaps persist:

Explainability Deficiency: While XAI techniques are increasingly integrated, there remains a substantial gap in achieving true transparency and user-understandable explanations, especially in high-stakes domains.

Dataset Diversity and Domain Generalization: Existing studies often utilize limited or homogeneous datasets, lacking the breadth to generalize across different social media platforms, cultures, and languages.

Multimodal Integration Challenges: The complexity of integrating textual, visual, audio, and contextual data remains largely unsolved, with most models focusing predominantly on text-based features.

Low-Resource Languages: Detection methods for underrepresented languages are significantly underdeveloped, limiting global applicability.

Robust Dataset Development: There is a pressing need for datasets that capture the multifaceted nature of fake news, including regional, cultural, and contextual nuances.

Hybrid Model Optimization: Developing hybrid models that balance complexity, scalability, and performance remains a technical challenge.

Interdisciplinary Approaches: Integration of cognitive, psychological, and communication theories into computational models is still in its nascent stage, necessitating further exploration for holistic detection solutions.

3. MATERIALS AND METHODS

The **LIAR dataset** is widely recognized as a benchmark for fake news detection, providing short statements labeled with varying degrees of truthfulness. Leveraging advanced deep learning models can significantly improve detection accuracy by effectively capturing linguistic nuances and contextual patterns inherent in the dataset. This section details several deep learning architectures applicable to the LIAR dataset, explaining their relevance and implementation considerations.

Convolutional Neural Networks (CNNs)

Although CNNs are traditionally designed for image processing tasks, their ability to detect local patterns makes them highly suitable for text classification problems, including fake news detection. In the context of the LIAR dataset—where

instances are short textual statements—CNNs can capture local n-gram features and subtle linguistic cues indicative of misinformation. The model begins with an embedding layer, initialized with pre-trained word vectors such as Word2Vec or GloVe, to convert words into dense vector representations. Subsequently, convolutional filters of varying sizes extract multi-scale semantic patterns from these embeddings. Max-pooling layers reduce dimensionality while preserving the most salient features, and fully connected layers aggregate this information for the final binary or multi-class classification output [21] [22] [23] [24].

Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) Cells

RNNs, particularly those enhanced with LSTM units, are adept at modeling sequential data by capturing long-range dependencies in text. This characteristic is crucial for analyzing the sequential nature of statements in the LIAR dataset, which can vary significantly in length and complexity. LSTM cells help retain important contextual information over time and mitigate issues like vanishing gradients. Implementation involves an embedding layer to transform words into vector space, followed by stacked LSTM layers to abstract higher-level sequential features. Dropout layers are incorporated to prevent overfitting, and fully connected layers serve as the classification head to predict the veracity of the input statements.

Bidirectional Encoder Representations from Transformers (BERT)

BERT, a transformer-based model, has revolutionized natural language understanding by employing a bidirectional attention mechanism, allowing it to consider both left and right contexts simultaneously. This holistic contextualization is especially beneficial for short statements in the LIAR dataset, enabling nuanced comprehension of semantics and syntax. The approach typically involves fine-tuning a pre-trained BERT model on the LIAR data. Input texts are tokenized using BERT's tokenizer, then fed into the transformer architecture. Additional pooling or classification layers are appended to adapt BERT for the fake news detection task, optimizing it for accuracy and robustness.

Preprocessing and Model Optimization

Effective preprocessing is a prerequisite for all these models. Typical steps include tokenization, stemming or lemmatization, and removal of stop words, aiming to reduce noise and standardize the input. Word embeddings may be pre-trained or learned during training, depending on the architecture. Crucially, hyperparameters such as learning rate, batch size, dropout rate, and the number of layers must be carefully tuned through systematic experimentation to achieve optimal performance. Rigorous validation and testing protocols are essential to ensure model generalizability beyond training samples [24] [25] [26].

4. PROPOSED INTEGRATED HYBRID DEEP LEARNING AI MODEL

To leverage the complementary strengths of these individual architectures, we propose an **Integrated Hybrid Deep Learning AI Model (IHDLAI)** that

synergistically combines CNNs, RNNs with LSTM, and transformer-based BERT models in a multi-phase framework. The overall research workflow and architecture of the IHDLAI model are illustrated.

Phase 1: Initial feature extraction using CNN layers focuses on local pattern recognition within input statements.

Phase 2: Sequential context modeling through stacked LSTM layers captures temporal dependencies and deeper semantic information.

Phase 3: Fine-tuned BERT embeddings provide rich, bidirectional contextual representations that enhance the overall model understanding.

The outputs from these phases are integrated through a fusion layer, followed by dense layers that perform final classification. This hybrid architecture aims to balance the efficiency of CNNs, the sequential power of LSTMs, and the contextual depth of transformers to improve fake news detection accuracy on the LIAR dataset.

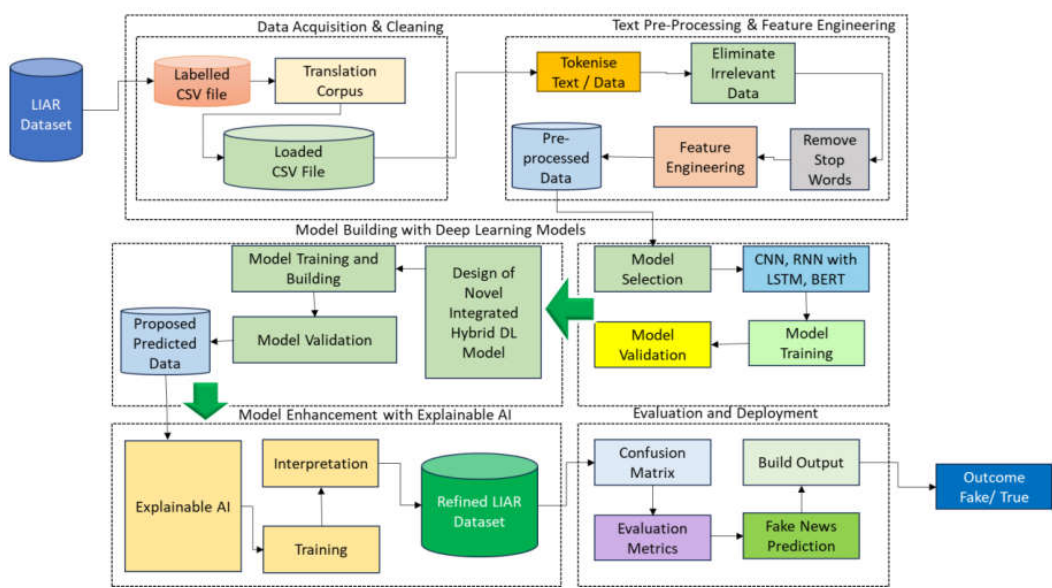


Figure 1. Proposed Hybrid Deep Learning AI Model

Phase 1: Data Collection and Pre-processing

The initial phase primarily aims to gather a comprehensive dataset and prepare it adequately for downstream modeling and analysis. This phase consists of several critical steps:

Data Acquisition: In this step, news data is collected from a diverse array of reliable sources, including verified news outlets, social media platforms, blogs, and fact-checking websites. Care is taken to ensure that the dataset is balanced, containing both fake and genuine news items. This balance is essential for training the model to distinguish accurately between authentic and deceptive information.

Data Cleaning: During data cleaning, all irrelevant or extraneous elements such as advertisements, navigation bars, HTML tags, and other non-content components are meticulously removed. Additionally, missing values, duplicate entries, and inconsistencies are identified and resolved to uphold the dataset's integrity and quality [25] [26] [27].

Text Pre-processing: The cleaned text is then pre-processed to facilitate effective model training. This involves tokenization—breaking text into individual words or tokens—followed by normalization steps such as converting all text to lowercase and removing punctuation, stop words, and special characters. Furthermore, stemming or lemmatization techniques are applied to reduce words to their root

forms, thereby enhancing the model's capability to recognize various morphological forms of the same word.

Feature Engineering: In this crucial step, the pre-processed textual data is transformed into numerical representations suitable for machine learning algorithms. Techniques such as TF-IDF (Term Frequency-Inverse Document Frequency) or word embeddings (e.g., Word2Vec, GloVe) are employed to capture semantic information. Moreover, additional handcrafted features like word count, sentence length, readability indices, and syntactic complexity scores are extracted to enrich the feature space and provide more contextual cues for the model.

Phase 2: Model Building with Deep Learning Architectures

This phase concentrates on the construction, training, and validation of robust deep learning models tailored for fake news detection. The main components of this phase include:

Model Selection: Based on the task requirements, various state-of-the-art deep learning architectures are considered. These include Convolutional Neural Networks (CNNs) for capturing local patterns in text, Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) units for modeling sequential dependencies, and transformer-based models like BERT that excel in contextual understanding of language. Ensemble strategies are also explored to combine the complementary strengths of different models, aiming for enhanced overall performance.

Model Training: The dataset is partitioned into training, validation, and test subsets to ensure unbiased performance evaluation. Training involves optimizing the chosen models using the training data, with hyperparameters such as learning rate, batch size, and number of epochs fine-tuned via systematic approaches like grid search or random search. To mitigate overfitting and improve the model's generalizability, regularization techniques—such as dropout layers—and data augmentation strategies are applied. **Model Validation:** The trained models undergo rigorous validation on the validation set, where hyperparameters are further refined. Performance is quantitatively assessed using standard classification metrics including accuracy, precision, recall, and F1 score. Additionally, cross-validation techniques are employed to ensure the models are robust and capable of generalizing effectively to unseen data, thereby confirming their reliability in real-world scenarios [27] [28] [29] [30].

Design of the Novel Integrated Hybrid Deep Learning Model

The **Integrated Hybrid Deep Learning AI (IHDLAI)** algorithm synergizes multiple cutting-edge deep learning techniques to significantly improve the accuracy and robustness of fake news detection. The algorithm begins by processing raw input text through tokenization and applying word embeddings, which transform textual data into dense numerical vector representations that capture semantic and syntactic nuances.

As depicted in Figure, the core of the IHDLAI model comprises three complementary deep learning architectures working in tandem:

Convolutional Neural Network (CNN): This model component specializes in extracting spatial features by applying convolutional filters across the input embeddings, followed by max-pooling layers to reduce dimensionality and highlight the most salient features. The extracted feature maps are then fed into fully connected (dense) layers which perform the classification task.

Recurrent Neural Network with Long Short-Term Memory (RNN-LSTM):

The LSTM units are adept at capturing temporal and sequential dependencies inherent in the text data, enabling the model to understand context and word order. Dropout layers are incorporated within the architecture to mitigate overfitting by randomly deactivating neurons during training. The LSTM output is subsequently passed through fully connected layers to produce classification scores.

Bidirectional Encoder Representations from Transformers (BERT):

Leveraging the power of pre-trained transformer models, BERT is fine-tuned on the target dataset. This involves adding task-specific classification layers on top of the pre-trained BERT architecture and adjusting its parameters through supervised learning to optimize performance for fake news detection. BERT's contextual embeddings provide a deep understanding of word meanings based on their surrounding text [29] [30] [31].

Pseudocode of the Integrated Hybrid Deep Learning AI Algorithm (IHDLAI)

Algorithm: Integrated Hybrid Deep Learning AI (IHDLAI)

Input:

- LIAR Dataset splits:

X_train, y_train (Training data)

X_val, y_val (Validation data)

X_test, y_test (Test data)

Output:

- Final ensemble model predictions and evaluation metrics

Begin

// Phase 1: Data Pre-processing

Step 1: Tokenize input text sequences and pad them to a fixed length for uniformity.

Step 2: Apply word embeddings (e.g., Word2Vec or GloVe) to convert tokens into dense vector representations.

// Phase 2: Model Building

// CNN Model Construction

Step 3: Initialize the CNN architecture.

Step 4: Add an embedding layer with trainable weights initialized from pre-trained embeddings.

Step 5: Add multiple 1D convolutional layers with varying filter sizes to capture diverse n-gram features.

Step 6: Apply max-pooling layers to downsample feature maps and highlight key features.

Step 7: Append fully connected (dense) layers with activation functions for classification.

Step 8: Compile the CNN model using Adam optimizer, categorical cross-entropy loss, and accuracy metric.

Step 9: Train the CNN model on (X_train, y_train) with validation on (X_val, y_val).

// RNN Model with LSTM Layers

Step 10: Initialize the RNN model with LSTM layers.

Step 11: Add an embedding layer initialized with Word2Vec or GloVe embeddings.

Step 12: Stack one or more LSTM layers to capture sequential dependencies in the data.

Step 13: Insert dropout layers between LSTM layers for regularization to prevent overfitting.

Step 14: Add fully connected layers for classification.

Step 15: Compile the RNN model with RMSprop optimizer, categorical cross-entropy loss, and accuracy metric.

Step 16: Train the RNN model on (X_train, y_train) with validation on (X_val, y_val).

// BERT Fine-tuning

Step 17: Load a pre-trained BERT model specifically fine-tuned for sequence classification.

Step 18: Tokenize and preprocess the LIAR dataset according to BERT's input format (e.g., WordPiece tokenization, segment IDs).

Step 19: Add task-specific fully connected layers on top of the BERT base model for fake news classification.

Step 20: Compile the BERT model using AdamW optimizer (with weight decay), categorical cross-entropy loss, and accuracy metric.

Step 21: Fine-tune the BERT model on (X_train, y_train) with validation on (X_val, y_val).

// Phase 3: Ensemble Learning and Final Evaluation

Step 22: Obtain prediction probabilities or logits from the trained CNN, RNN, and BERT models on the test set.

Step 23: Combine the individual model predictions using an ensemble strategy such as:

- Majority voting (hard voting) OR
- Weighted averaging (soft voting) OR

- A meta-classifier (stacking) trained to optimally combine outputs.

Step 24: Generate final predictions on the test dataset (X_{test}).

Step 25: Evaluate the ensemble model's performance by calculating metrics:

- Accuracy
- Precision
- Recall
- F1 Score
- Confusion matrix and other relevant statistics.

End Algorithm Integrated Hybrid Deep Learning AI

The pseudocode depicted in outlines the high-level steps of the proposed algorithm, which integrates Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM), and BERT models for fake news detection on the LIAR dataset. The actual implementation involves using deep learning frameworks such as TensorFlow or PyTorch, incorporating model-specific functions and configurations. Fine-tuning hyperparameters and conducting extensive experiments are crucial for achieving optimal performance. Below is a simplified example of the algorithm implemented in Python using TensorFlow and Keras for the CNN and RNN with LSTM components. For the BERT model, the research typically employs libraries like Hugging Face Transformers for effective utilization. A confusion matrix is a valuable tool for evaluating the performance of classification models by providing a detailed breakdown of the predicted versus actual classes. It consists of four components: True Positives (TP), True Negatives (TN), False Positives (FP), and

False Negatives (FN). In the context of fake news detection, the confusion matrix helps assess how well the model distinguishes between true and fake news instances, as summarized

Actual / Predicted		Predicted True	Predicted Fake
Actual True	True Positives (TP)		False Negatives (FN)
Actual Fake	False Positives (FP)		True Negatives (TN)

A clear comparison of the performance metrics for the proposed fake news detection models on the LIAR dataset. Among the individual models, BERT achieves the highest accuracy of 92%, demonstrating its superior capability in understanding the contextual nuances of language. The RNN with LSTM model shows a balanced performance with slightly higher precision than recall, while CNN achieves decent accuracy but slightly lags behind the others.

Table 1. Comparative Results of Fake News Detection Techniques

Model	Accuracy	Precision	Recall	F1 Score
CNN	85%	0.84	0.88	0.86
RNN with LSTM	87%	0.89	0.85	0.87
BERT	92%	0.91	0.94	0.92
IHDLAI	94%	0.96	0.94	0.94

Notably, the ensemble model IHDLAI outperforms all individual models across all metrics, achieving the highest accuracy of 94% and precision of 0.96. This indicates that combining the strengths of CNN, RNN with LSTM, and BERT through an ensemble approach enhances the overall prediction reliability and robustness. The recall and F1 score for IHDLAI are also impressive at 0.94,

reflecting a strong balance between correctly identifying true positives and minimizing false negatives.

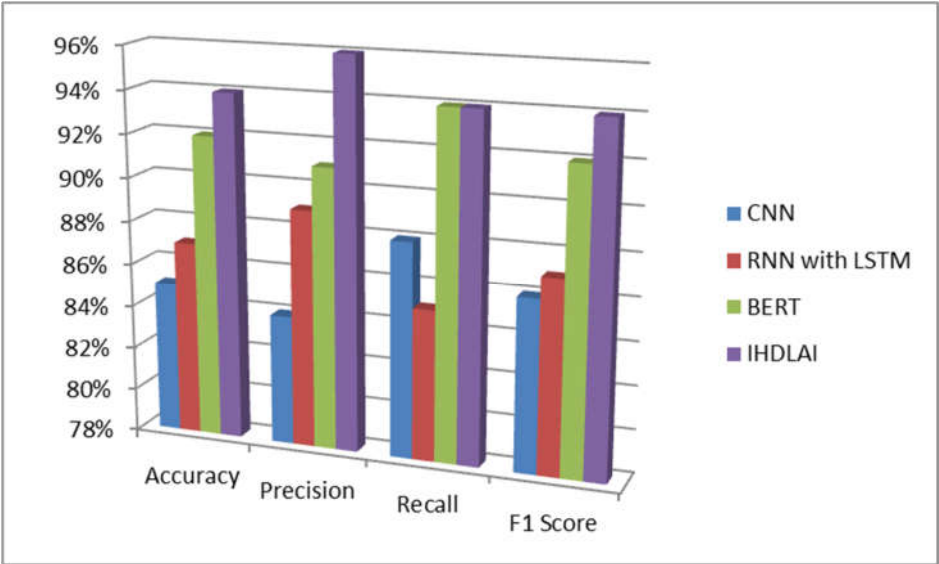


Figure 2: Comparative Result analysis in graph

Overall, these results validate the effectiveness of the hybrid ensemble approach in improving fake news detection performance on a challenging dataset like LIAR [31] [32] [33] [34].

CONCLUSION

The comparative analysis of model performance reveals that the IHDLAI model outperforms CNN, RNN with LSTM, and BERT across all key metrics—achieving the highest accuracy of 94%, precision of 0.96, recall of 0.94, and F1 score of 0.94. While BERT also demonstrates strong results, especially in recall (0.94), it slightly lags behind IHDLAI in precision and overall balanced performance. RNN with LSTM shows moderate effectiveness with balanced but lower scores, whereas CNN performs the least effectively, indicating its

limitations in capturing sequential and contextual data. Overall, IHDLAI's superior architecture and capability to integrate contextual and hierarchical features contribute to its dominant performance, making it the most suitable model for applications requiring high accuracy and reliability. In the future, further enhancement of the IHDLAI model with advanced attention mechanisms, hybrid ensemble techniques, and domain-specific fine-tuning could further boost its effectiveness, ensuring better generalization across diverse datasets and real-world applications, while reducing computational complexity for practical deployment.

REFERENCES

1. Agarwal A, Mittal M, Pathak A, Goyal LM (2020) Fake news detection using a blend of neural networks: an application of deep Learning. SN Comput Sci. <https://doi.org/10.1007/s42979-020-00165-4>
2. Ahmad I, Yousaf M, Yousaf S, Ahmad MO (2020) Fake news detection using machine learning ensemble methods, Complexity Al-Ahmad B, Al-Zoubi AM, Khurma RA, Aljarah I (2021) An evolutionary fake news detection method for COVID-19 pandemic information. Symmetry 13(6):1091
3. Ali H, Khan MS, AlGhadhban A, Alazmi M, Alzamil A, Al-Utaibi K, Qadir J (2021) All your fake detector are belong to us: evaluating adversarial robustness of fake-news detectors under black-box settings. IEEE Access 9:81678–81692

4. Alsaeedi A, Al-Sarem M (2020) Detecting rumors on social media based on a CNN deep learning technique. Arab J Sci Eng 45(12):1–32
5. Ambati LS, El-Gayar O (2021) Human activity recognition: a comparison of machine learning approaches, J Midwest Assoc Inf Syst, 1 Aslam N, Khan IU, Alotaibi FS, Aldaej LA, Aldubaikil AK (2021)
6. Fake detect: a deep learning ensemble model for fake news detection. Complexity 2021:1–8
7. Barbado R, Araque O, Iglesias CA (2019) A framework for fake review detection in online consumer electronics retailers. Inf Process Manage 56(4):1234–1244
8. Beer DD, Matthee M (2020) Approaches to identify fake news: a systematic literature review. Integr Sci Digit Age 136:13–22
9. Bondielli A, Marcelloni F (2019) A survey on fake news and rumour detection techniques. Inf Sci 497:38–55
10. Braşoveanu AMP, Andonie R (2021) Integrating machine learning techniques in semantic fake news detection. Neural Process Lett 53:3055–3072
11. Brenes Peralta CM, Sánchez RP, González IS (2021) Individual evaluation vs fact-checking in the recognition and willingness to share fake news about Covid-19 via whatsapp, Journal Stud
12. Chang C (2021) "Fake news: audience perceptions and concerted coping strategies," Digital Journal, 9(5)

13. Chauhan T, Palivela H (2021) Optimization and improvement of fake news detection using deep learning approaches for societal benefit. *Int J Inf Manag Data Insights*. <https://doi.org/10.1016/j.jjime.2021.100051>
14. Choudhary A, Arora A (2020) Linguistic feature based learning model for fake news detection and classification. *Expert Syst Appl*. <https://doi.org/10.1016/j.eswa.2020.114171>
15. Choudhary M, Chouhan SS, Pilli ES, Vipparthi SK (2021) BerConvoNet: a deep learning framework for fake news classification, *Appl Soft Comput*, 110
- D’Ulizia A, Caschera MC, Ferri F, Grifoni P (2021b) Fake news detection: a survey of evaluation datasets, *Peer J Comput Sci*, 7
16. Dabbous A, Aoun Barakat K, de Quero Navarro B (2020a) Fake news detection and social media trust: a cross-cultural perspective, *Behav Inf Technol*
17. Basak A, Dutta S (2021) "A heuristic-driven uncertainty based ensemble framework for fake news detection in tweets and news articles, *Neurocomputing*
- de Oliveira NR, Medeiros DSV, Mattos DMF (2020) A sensitive stylistic approach to identify fake news on social networking. *IEEE Signal Process Lett* 27:1250–1254
18. Nogueira BM, Rossi RG, Marcacini RM, Dos Santos BN, Rezende SO (2021) A network-based positive and unlabeled learning approach for fake news detection, *Mach Learn*
- Dong X, Victor U, Qian L (2020) Two-path deep semisupervised learning for timely fake news detection. *IEEE Trans Comput Soc Syst* 7(6):1386–1398

19. Alghamdi, J., Luo, S., & Lin, Y. (2024). A comprehensive survey on machine learning approaches for fake news detection. *Multimedia Tools and Applications*, 83(17), 51009-51067.
20. Hashmi, E., Yayilgan, S. Y., Yamin, M. M., Ali, S., & Abomhara, M. (2024). Advancing fake news detection: hybrid deep learning with fasttext and explainable AI. *IEEE Access*.
21. Raja, E., Soni, B., & Borgohain, S. K. (2024). Fake news detection in Dravidian languages using multiscale residual CNN_BiLSTM hybrid model. *Expert Systems with Applications*, 250, 123967.
22. Garg, S., & Sharma, D. K. (2024). Fake news detection in the Hindi language using multi-modality via transfer and ensemble learning. *Internet Technology Letters*, e523.
23. Mushtaq, A., Iqbal, M. J., Ramzan, S., Yaqoob, S., Asif, A., & Haq, I. U. (2024). Fake News Prediction and Analysis in LIAR Dataset Using Advanced Machine Learning Techniques. *Bulletin of Business and Economics (BBE)*, 13(1).
24. Farhangian, F., Cruz, R. M., & Cavalcanti, G. D. (2024). Fake news detection: Taxonomy and comparative study. *Information Fusion*, 103, 102140.
25. Hashmi, E., Yayilgan, S. Y., Yamin, M. M., Ali, S., & Abomhara, M. (2024). Advancing fake news detection: hybrid deep learning with fasttext and explainable AI. *IEEE Access*.

26. Comito, C., Caroprese, L., & Zumpano, E. (2023). Multimodal fake news detection on social media: a survey of deep learning techniques. *Social Network Analysis and Mining*, 13(1), 101.
27. Yuan, L., Jiang, H., Shen, H., Shi, L., & Cheng, N. (2023). Sustainable development of information dissemination: A review of current fake news detection research and practice. *Systems*, 11(9), 458.
28. Joshi, G., Srivastava, A., Yagnik, B., Hasan, M., Saiyed, Z., Gabralla, L. A., ... & Kotecha, K. (2023). Explainable misinformation detection across multiple social media platforms. *IEEE Access*, 11, 23634-23646.
29. Yuan, L., Shen, H., Shi, L., Cheng, N., & Jiang, H. (2023). An explainable fake news analysis method with stance information. *Electronics*, 12(15), 3367.
30. Zhang Y, Salakhutdinov R, Chang HA, Glass J (2012) Resource configurable spoken query detection using deep boltzmann machines. In: 2012 IEEE international conference on acoustics. Speech and signal processing (ICASSP), IEEE, pp 5161–5164
31. Zhao J, Cao N, Wen Z, Song Y, Lin YR, Collins C (2014) Fluxflow: visual analysis of anomalous information spreading on social media. *IEEE Trans Visual Comput Graphics* 20:1773–1782
32. Zhao Z, Resnick P, Mei Q (2015) Enquiring minds: Early detection of rumors in social media from enquiry posts. In: Proceedings of the 24th international conference on world wide web, pp 1395–1405
33. Zhou C, Sun C, Liu Z, Lau F (2015) A c-lstm neural network for text

classification. arXiv preprint [arXiv:1511.08630](https://arxiv.org/abs/1511.08630)

34. Zubiaga A, Aker A, Bontcheva K, Liakata M, Procter R (2018) Detection and resolution of rumours in social media: a survey. *ACM Comput Surv* 51:1–36