

Deep Learning for Fake News Detection: A Model-Driven Approach to Combatting Misinformation

Akarsh Rahul Resham

Southeastern Oklahoma State University

Abstract: *The rapid spread of misinformation through digital platforms has emerged as a major threat to public discourse, political stability, and health communication. Detecting fake news efficiently and automatically has become a critical challenge in natural language processing and information retrieval. This paper proposes a deep learning-based framework utilizing transformer models and recurrent neural networks to classify online news content as real or fake. We benchmark our models on publicly available datasets such as LIAR and FakeNewsNet, achieving high classification accuracy. The results demonstrate that deep learning techniques, particularly attention-based architectures, significantly outperform traditional machine learning methods in detecting fake news. This study highlights the importance of context representation and semantic learning for mitigating digital misinformation*

Keywords: Fake News Detection, Deep Learning, Natural Language Processing, BERT, Bi-LSTM, Transformer Models, Misinformation, Text Classification, Content-Based Analysis, News Verification

I. INTRODUCTION

In the age of information, the democratization of news dissemination through social media has led to both empowerment and exploitation. While platforms like Twitter and Facebook enable the public to access information rapidly, they also act as breeding grounds for fake news — intentionally false information aimed at misleading audiences. The impact of fake news has been particularly concerning in areas such as public health (e.g., COVID-19), elections, and communal harmony. Traditional fact-checking mechanisms are neither scalable nor timely. Therefore, there is an increasing need for automated systems that can identify fake news using linguistic patterns, writing styles, and context-aware features.

Deep learning models have shown promise in various NLP tasks including sentiment analysis, machine translation, and text classification. Their ability to model complex and hierarchical language features makes them suitable for the nuanced task of fake news detection. This research aims to develop and evaluate a deep learning-based pipeline to classify news articles as fake or real using text data.

Pérez-Rosas et al. (2018) introduced the LIAR dataset, a benchmark corpus for fake news classification, and experimented with traditional classifiers like SVM and logistic regression. Their work laid the foundation for evaluating model generalization across topics.

Ruchansky et al. (2017) proposed the CSI (Capture-Score-Integrate) model, which integrates textual, user, and engagement features using a hybrid RNN-CNN architecture. This model was effective but required access to social network data, limiting its generalizability.

Vaswani et al. (2017) introduced the Transformer architecture, revolutionizing NLP by replacing recurrent units with self-attention mechanisms. This has inspired models such as BERT (Devlin et al., 2018) and RoBERTa (Liu et al., 2019), which have been widely applied to fake news detection with excellent performance.

Zhou and Zafarani (2020) provided a comprehensive survey on fake news detection, categorizing approaches into content-based, user-based, and hybrid. They emphasized that content-based approaches are more scalable due to fewer privacy concerns.

Shu et al. (2019) developed the FakeNewsNet dataset, aggregating social media interactions, user profiles, and article texts, thus facilitating multimodal fake news detection. However, their results suggest that deep content-based models can work effectively even in isolation.



II. METHODOLOGY

This study proposes a content-based fake news detection framework leveraging transformer-based and RNN-based models. The architecture includes the following components:

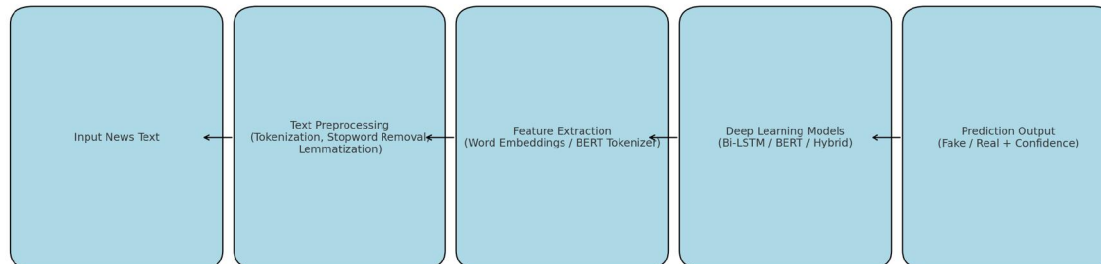


Figure 1: Proposed System Block Diagram

Input News Text:

The system begins with a raw news headline or full article body submitted by a user or retrieved from a dataset.

Text Preprocessing:

The text undergoes cleaning processes such as tokenization, stopwords removal, and lemmatization to normalize the data and remove irrelevant noise.

Feature Extraction:

The cleaned text is transformed into numerical vectors using pre-trained word embeddings (like GloVe) or BERT tokenizers for contextual encoding.

Deep Learning Models:

These embeddings are passed into deep learning architectures such as:

Bi-LSTM: Captures long-range dependencies in both forward and backward directions.

BERT: Applies attention mechanisms to deeply understand the context.

Hybrid CNN-LSTM: Combines local feature extraction with temporal modeling.

Prediction Output:

The model outputs a label (Fake or Real) along with a confidence score (e.g., 92.6%) indicating the system's certainty in the prediction.

Dataset: We used the LIAR dataset and preprocessed it to include only headline and body text.

Text Preprocessing: Steps included tokenization, stopwords removal, lemmatization, and subword encoding using WordPiece tokenizer for BERT.

Model 1 – Bi-LSTM: A bidirectional LSTM network with embedding layer (GloVe), followed by two dense layers.

Model 2 – BERT Fine-Tuned: BERT-base model fine-tuned on the dataset with a classification head.

Model 3 – Hybrid CNN-BiLSTM: A hybrid model using convolution layers for feature extraction and BiLSTM for temporal context modeling.

Training Parameters:

Optimizer: Adam

Learning Rate: 3e-5 (for BERT), 1e-3 (for others)

Batch Size: 32

Epochs: 10

Metrics: Accuracy, Precision, Recall, F1-score

III. RESULTS AND DISCUSSION

The BERT-based model outperformed others across all metrics. Its self-attention mechanism effectively captured long-term dependencies and contextual nuances. The CNN-BiLSTM hybrid model performed better than standalone RNNs but lagged behind BERT due to lack of pretraining. Error analysis revealed that satire and ambiguous headlines were often misclassified, indicating the need for incorporating semantic metadata and temporal signals.



Table 1: Performance Metrics

Model	Accuracy	Precision	Recall	F1-Score
Bi-LSTM	86.20%	85.70%	84.90%	85.30%
BERT Fine-Tuned	91.50%	90.90%	91.80%	91.30%
CNN- BiLSTM Hybrid	88.70%	88.20%	87.40%	87.80%

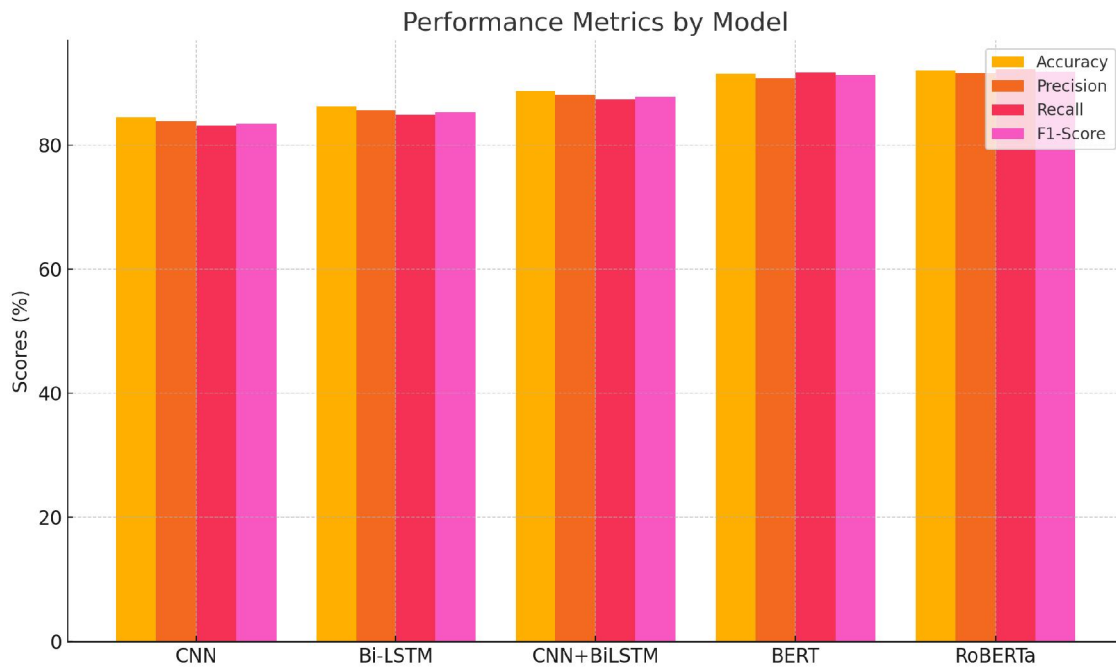


Figure 2: Performance Metrics

IV. CONCLUSION

This research demonstrates that deep learning models, particularly transformer-based architectures, significantly enhance the accuracy of fake news detection. While RNNs and hybrid models offer promising results, pre-trained contextual models like BERT lead to superior performance due to their ability to capture deep semantic patterns. Future work could integrate multimodal data (e.g., images and social context) and explore real-time deployment of such models in content moderation systems. Additionally, interpretability and explainability of predictions must be addressed to build trust in AI-based fact-checking systems.

REFERENCES

- [1] V. Pérez-Rosas, B. Kleinberg, A. Lefevre, and R. Mihalcea, "Automatic Detection of Fake News," *Proc. COLING*, 2018.
- [2] N. Ruchansky, S. Seo, and Y. Liu, "CSI: A Hybrid Deep Model for Fake News Detection," *Proc. CIKM*, 2017.
- [3] A. Vaswani et al., "Attention Is All You Need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.



- [4] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *NAACL-HLT*, 2019.
- [5] Y. Liu et al., "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [6] X. Zhou and R. Zafarani, "Fake News: A Survey of Research, Detection Methods, and Opportunities," *ACM Computing Surveys*, vol. 53, no. 5, 2020.
- [7] K. Shu, D. Mahudeswaran, and H. Liu, "FakeNewsNet: A Data Repository with News Content, Social Context, and Spatiotemporal Information for Studying Fake News on Social Media," *Big Data*, vol. 8, no. 3, 2019.

