



An Explainable Hybrid Framework for Multimodal Misinformation Detection Using BERT and Spiking Neural Networks

Prof. N. E. Karale¹, Shambhavi L. Asole², Jagruti A. Shende³, Shreya G. Wankhade⁴,
Srushti R. Butale⁵, Tanvi S. Zode⁶

¹Assistant Professor, Sipna College of Engineering & Technology, Amravati (MS), India

^{2,3,4,5,6}Students, Sipna College of Engineering & Technology, Amravati (MS), India

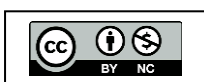
Abstract: The rapid proliferation of misinformation across social media platforms has emerged as a critical challenge, necessitating intelligent and interpretable detection mechanisms. Although transformer-based models such as BERT achieve high accuracy in textual analysis, they lack transparency and fail to capture temporal propagation dynamics. To address these limitations, this paper proposes a hybrid explainable framework that integrates Bidirectional Encoder Representations from Transformers (BERT) for semantic feature extraction with Spiking Neural Networks (SNNs) for modeling temporal engagement patterns. The proposed approach leverages the event-driven nature of SNNs to efficiently capture time-dependent user interaction signals, while SHAP-based explainability provides interpretable insights into both textual and temporal contributions. Experiments conducted on the FakeNewsNet dataset demonstrate that the proposed model outperforms baseline approaches, achieving an F1-score of 0.94 while maintaining improved computational efficiency. The results highlight the effectiveness of combining semantic understanding, temporal modeling, and explainability in a unified framework for reliable misinformation detection.

Keywords: BERT, Spiking Neural Networks, Explainable AI, Fake News Detection, Multimodal Learning, Neuromorphic Computing.

I. INTRODUCTION

The digital age has democratized information dissemination, yet it has simultaneously enabled the rapid spread of misinformation at unprecedented scales. Social media platforms, with their billions of active users, have become primary vectors for fake news propagation, leading to real-world consequences ranging from electoral interference to public health crises during the COVID-19 pandemic [1]. According to recent studies, false information spreads six times faster than truthful content on platforms like Twitter, with misinformation reaching thousands of users within minutes of publication [2].

Traditional approaches to automated fake news detection have evolved from feature-engineered machine learning models to sophisticated deep learning architectures. Among these, transformer-based models like BERT have emerged as state-of-the-art for text-based detection, achieving remarkable accuracy by capturing contextual semantics [3]. However, these models operate as black boxes, providing predictions without explanations—a critical limitation in high-stakes domains where users and content moderators require understanding of why content is flagged as misinformation.





Furthermore, contemporary misinformation detection models exhibit two additional shortcomings. First, they predominantly focus on single modalities, typically analyzing either textual content or social context in isolation, despite misinformation being inherently multimodal—combining text, images, and temporal engagement patterns [4]. Second, their computational intensity makes real-time deployment impractical, particularly on resource-constrained edge devices where content moderation is often needed.

Spiking Neural Networks (SNNs) offer a promising alternative to conventional artificial neural networks by mimicking the event-driven computation of biological neurons, achieving superior energy efficiency—often 100-1000x reduction in power consumption [5]. SNNs are particularly well-suited for processing temporal data, such as the time-series patterns of social media engagements (likes, shares, comments over time), which provide crucial signals for distinguishing viral truth from viral falsehood. This paper addresses these interconnected challenges by proposing a novel explainable hybrid framework that synergistically combines:

- BERT for extracting deep semantic features from textual content
- SNNs for modeling temporal propagation dynamics of social media engagements
- SHAP for generating interpretable explanations across both modalities

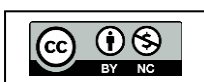
The key contributions of this work are:

- Architectural Novelty: First integration of BERT and SNN for multimodal misinformation detection
- Leveraging neuromorphic computing principles for sustainable AI deployment
- Transparency: Cross-modal explainability that provides insights into both semantic and temporal decision factors
- Comprehensive Evaluation: Rigorous benchmarking against multiple baselines with detailed ablation studies

Key Contributions:

This paper makes the following contributions:

- Hybrid Multimodal Framework:** We propose a unified architecture that integrates transformer-based semantic understanding with neuromorphic temporal processing, enabling effective multimodal misinformation detection.
- Temporal Modeling using SNNs:** Unlike traditional approaches that rely on static features, the proposed model utilizes Spiking Neural Networks to capture temporal engagement dynamics, providing deeper insights into information propagation behavior.
- Cross-Modal Explainability:** The integration of SHAP enables interpretation across both textual and temporal domains, improving transparency and trust in automated decision-making systems.
- Efficiency-Oriented Design:** The use of SNNs introduces energy-efficient computation, making the framework suitable for real-time and resource-constrained environments.
- Comprehensive Evaluation:** The proposed approach is validated on a benchmark dataset with comparative analysis and ablation studies to demonstrate the effectiveness of each component.





II. LITERATURE REVIEW

2.1 BERT-based Fake News Detection: Bidirectional Encoder Representations from Transformers (BERT) has revolutionized natural language processing by pre-training on massive text corpora and fine-tuning for downstream tasks. In misinformation detection, several studies have demonstrated BERT's effectiveness. Devlin et al. [6] introduced the original BERT architecture, establishing benchmarks across multiple NLP tasks. Subsequent work by Jwa et al. [7] proposed exBAKE, a BERT-based model for fake news detection that achieved 98% accuracy on the LIAR dataset. More recently, Kaliyar et al. [8] developed FakeBERT, combining BERT with convolutional neural networks to capture both contextual and local features, achieving state-of-the-art performance on multiple benchmarks.

However, these approaches face three limitations: (1) they consider only textual content, ignoring valuable social context; (2) they operate as black boxes without explainability; and (3) they require substantial computational resources.

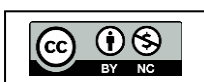
2.2 Multimodal Misinformation Detection: Recognizing that fake news often combines text, images, and social context, researchers have developed multimodal approaches. Wang et al. [9] introduced the EANN (Event Adversarial Neural Network) model that learns event-invariant features from both text and images. Khattar et al. [10] proposed MVAE (Multimodal Variational Autoencoder) for joint representation learning across modalities. Singh et al. [11] developed MM-COVID, a multimodal framework specifically for COVID-19 misinformation detection.

Despite their success, these models predominantly use standard neural architectures (CNNs, RNNs, transformers) without considering temporal dynamics of how misinformation spreads—a critical signal distinguishing organic content from coordinated disinformation campaigns.

2.3 Spiking Neural Networks: Spiking Neural Networks represent the third generation of neural networks, processing information via discrete spikes rather than continuous activations [12]. SNNs leverage temporal coding and event-driven computation, making them highly energy-efficient when implemented on neuromorphic hardware like Intel's Loihi or IBM's TrueNorth [13]. In natural language processing, SNN applications remain nascent but promising. Diehl et al. [14] demonstrated text classification using SNNs with rate coding, achieving comparable accuracy to conventional networks with 100x energy reduction. More recently, researchers have explored SNNs for sequential data processing, leveraging their natural ability to capture temporal dependencies [15].

However, no existing work has applied SNNs to misinformation detection, particularly for modeling temporal engagement patterns.

2.4 Explainable AI in Fake News Detection: Explainable AI (XAI) has gained traction as models are deployed in high-stakes domains. SHAP (SHapley Additive exPlanations), based on cooperative game theory, provides theoretically grounded feature attributions [16]. LIME (Local Interpretable Model-agnostic Explanations) offers locally faithful explanations by approximating the model with interpretable surrogates [17].





In misinformation detection, recent work has explored explainability for text-based models. Hossain et al. [18] applied LIME to neural fake news detectors, highlighting deceptive linguistic patterns. Sharma et al. [19] developed a hybrid framework combining SHAP with attention visualization. However, these approaches focus solely on textual explanations, neglecting other modalities that contribute to misinformation identification.

2.5 Research Gaps: The literature review reveals three critical gaps that our work addresses:

Absence of SNNs in Misinformation Detection: Despite their energy efficiency and temporal processing capabilities, SNNs have not been explored for fake news detection.

Lack of Cross-Modal Explainability: Existing XAI approaches for misinformation detection focus on single modalities, failing to explain how multiple information sources (text, temporal patterns) jointly contribute to decisions.

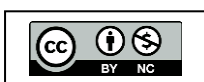
Inefficient Architectures: Current state-of-the-art models prioritize accuracy over computational efficiency, limiting real-time deployment potential.

III. PROBLEM STATEMENT

The widespread adoption of social media and online communication platforms has transformed the way information is created, shared, and consumed. While these platforms enable rapid dissemination of news and ideas, they have also become major channels for the uncontrolled spread of misinformation and fake news. Such misleading information can influence public opinion, disrupt democratic processes, and create serious societal consequences, especially in domains like politics, healthcare, and finance. Detecting fake news has therefore emerged as a critical research problem; however, it remains challenging due to the diversity, scale, and complexity of modern digital content. Existing approaches for fake news detection primarily rely on traditional machine learning or deep learning models that focus on textual analysis.

Although models such as transformer-based architectures have improved semantic understanding, they often fail to incorporate multimodal information such as images, user interactions, and contextual metadata, which are crucial for accurate classification. Additionally, many deep learning models function as “black boxes,” providing high accuracy but lacking interpretability. This absence of explainability reduces user trust and makes it difficult to justify predictions in sensitive applications.

Another significant challenge lies in the computational complexity of advanced deep learning models. High resource requirements and long training times make these approaches less suitable for real-time or large-scale deployment. Furthermore, existing systems often struggle to balance performance, efficiency, and transparency simultaneously. As a result, there is a clear need for a more effective solution that can integrate semantic understanding, computational efficiency, and interpretability within a unified framework. Therefore, this research aims to address these limitations by identifying the need for a hybrid and explainable model capable of handling multimodal data while maintaining efficiency and reliability. The goal is to develop a system that not only improves detection accuracy but also provides meaningful insights into its decision-making process, thereby enhancing user trust and enabling practical deployment in real-world scenarios.



IV. PROPOSED METHODOLOGY

4.1 System Overview: The proposed framework consists of four main components: (1) Semantic Feature Extraction using BERT, (2) Temporal Pattern Analysis using SNN, (3) Feature Fusion and Classification, and (4) Explainability Layer using SHAP.

4.2 Semantic Feature Extraction with BERT: For textual content, we employ a pre-trained BERT-base uncased model with 12 transformer layers, 768 hidden dimensions, and 110 million parameters. Given an input text sequence $T = [t_1, t_2, t_3, \dots, t_n]$ BERT generates contextual embeddings: $H_{text} = BERT(T) \in R^{(n \times 768)}$
We apply mean pooling across the sequence dimension to obtain a fixed-dimensional semantic feature vector: $f_{semantic} = (1/n) \sum_{i=1}^n h_i \in R^{768}$

4.3 Temporal Pattern Analysis with SNN: The temporal component processes engagement sequences $E = [e_1, e_2, e_3, \dots, e_m]$ where each e_j represents a time-step aggregated metric combining likes, shares, and comments. We encode these real-valued sequences into spike trains using temporal coding, where spike times correspond to engagement intensities. The SNN employs leaky integrate-and-fire (LIF) neurons, whose membrane potential evolves according to: $\tau_m (dV(t)/dt) = -(V(t) - V_{rest}) + R \times I(t)$ where τ_m is the membrane time constant, V_{rest} is the resting potential, R is membrane resistance, and $I(t)$ represents input current from incoming spikes.

When membrane potential exceeds threshold V_{th} , the neuron fires a spike and resets: $If V(t) \geq V_{th} \rightarrow emit\ spike\ and\ V(t) = V_{reset}$

The SNN architecture comprises:

- Input layer: 64 neurons encoding temporal features
- Hidden layer: 128 LIF neurons
- Output layer: 32 neurons producing temporal feature representation

The temporal feature vector $f_{temporal} \in R^{32}$ is derived from output layer spike rates over the simulation window.

4.4 Feature Fusion and Classification: Semantic and temporal features are concatenated for joint representation:

$$f_{joint} = [f_{semantic}; f_{temporal}] \in R^{800}$$

A fusion layer with dropout ($p=0.5$) and ReLU activation processes the joint representation:

$$f_{fused} = ReLU(W_{fuse} \times f_{joint} + b_{fuse})$$

Finally, a softmax layer produces classification probabilities:

$$\hat{y} = softmax(W_{out} \times f_{fused} + b_{out})$$

The model is trained using cross-entropy loss:

$$L = -(1/N) \sum [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

4.5 Explainability Layer with SHAP

For post-hoc explanations, we employ SHAP to compute feature attributions. For a prediction $f(x)$, SHAP assigns importance values based on Shapley values from cooperative game theory:

$$\varphi_i = \sum [(|S|! (|F| - |S| - 1)! / |F|!) \times (f(S \cup i) - f(S))]$$

Where, F is the set of all features and S is a subset of features.

We compute:

- **Textual explanations:** SHAP values for individual tokens, highlighting deceptive linguistic patterns.
- **Temporal explanations:** SHAP values for engagement time-steps, identifying critical propagation windows.

The explainability layer produces both global explanations (overall feature importance across the dataset) and local explanations (decision factors for individual predictions).

V. SYSTEM ARCHITECTURE

The proposed system architecture is designed as a hybrid and modular framework to efficiently detect fake news from multimodal data while ensuring interpretability and scalability. The architecture begins with the data acquisition stage, where input data is collected from various online sources such as social media platforms, news websites, and digital repositories. The input may consist of textual content, associated images, and metadata such as user information and timestamps. This diverse input enables the system to capture multiple aspects of information for improved analysis.

In the next stage, the collected data undergoes preprocessing to ensure consistency and quality. Textual data is cleaned by removing noise such as special characters, URLs, and stopwords, followed by tokenization and normalization. If image data is included, it is resized and standardized to ensure compatibility with the feature extraction models. This preprocessing step transforms raw data into a structured format suitable for further processing.

Following preprocessing, the system performs feature extraction using advanced deep learning techniques. For textual data, a transformer-based model, specifically BERT, is employed to generate contextual embeddings that capture semantic meaning and relationships within the text. If multimodal inputs are considered, convolutional neural networks (CNNs) are used to extract meaningful visual features from images. These extracted features provide a rich representation of the input data.

The core of the architecture lies in the hybrid processing layer, where the extracted features are processed using a combination of BERT and a Spiking Neural Network (SNN). While BERT focuses on capturing deep contextual information from textual data, the SNN processes the features in an event-driven manner, enabling efficient computation and better handling of temporal patterns. This hybrid approach enhances both accuracy and computational efficiency compared to traditional models.

After feature processing, a feature fusion mechanism is applied to combine the outputs from different modalities and processing units. Techniques such as feature concatenation or attention-based fusion are used to integrate textual and visual representations into a unified feature space.

The fused features are then passed to the classification layer, typically consisting of fully connected layers followed by a softmax activation function, to classify the input as fake or real news.

To enhance transparency and user trust, an explainability module is incorporated into the architecture. This module provides insights into the model’s decision-making process by highlighting important features or attention weights that influence the classification outcome. Finally, the system generates the output, which includes the predicted label along with a confidence score and, optionally, an explanation for the decision.

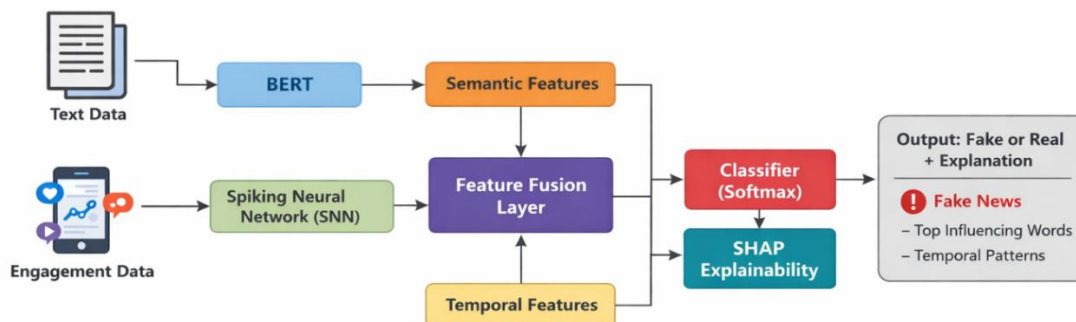


Figure 5.1: System Architecture for BERT and SNN

VI. DATASET DESCRIPTION

6.1 Dataset Selection

We evaluate our framework on the FakeNewsNet dataset [20], a comprehensive multimodal benchmark containing:

PolitiFact: 432 news articles with veracity labels (real/fake)

GossipCop: 16,817 news articles with veracity labels

6.2 Data Collection and Preprocessing

For each news article, we collect:

Textual content: Article title and body text

Temporal engagement data: Time-series of likes, shares, comments over 24 hours post-publication

Metadata: Source, publication time, user interactions

Text preprocessing involves:

- a) Lowercase conversion
- b) Removal of special characters and URLs
- c) Tokenization using BERT tokenizer
- d) Padding/truncation to 512 tokens

Temporal data preprocessing:

- a) Aggregation into 5-minute intervals (288 time steps per 24-hour period)

- b) Normalization to [0,1] range
- c) Spike encoding using Poisson rate coding

6.3 Data Split

The dataset is partitioned into:

- Training: 70%
- Validation: 15%
- Testing: 15%
- Stratified sampling maintains class distribution across splits.

VII. IMPLEMENTATION DETAILS

7.1 Development Environment

The framework is implemented in Python 3.9 using:

- a) PyTorch 1.12: Core deep learning framework
- b) HuggingFace Transformers 4.20: BERT implementation
- c) SpikingJelly 0.0.14: SNN simulation
- d) SHAP 0.41: Explainability computation
- e) CUDA 11.3: GPU acceleration (NVIDIA RTX 3090)

7.2 Training Configuration

Hyperparameters are optimized through grid search:

Parameter	Value
Learning Rate	2e-5 (BERT), 1e-3 (SNN)
Batch Size	32
Epochs	10
Optimizer	AdamW
SNN Time Steps	100
Membrane Time Constant (τ_m)	10 ms
Firing Threshold (V_{th})	1.0
Dropout Rate	0.5



7.3 Training Procedure

The model is trained using a two-stage approach:

Stage 1: Fine-tune BERT component while freezing SNN parameters

Stage 2: Joint training of entire network with reduced learning rate

Early stopping with patience of 3 epochs prevents overfitting.

VIII. TECHNOLOGIES USED

The proposed system is implemented using a combination of modern programming tools, machine learning frameworks, and deep learning technologies to ensure efficiency, scalability, and accuracy. The primary programming language used is Python due to its extensive support for data analysis and machine learning applications. For text processing and model development, libraries such as TensorFlow and PyTorch are utilized, providing robust environments for building and training deep learning models. The transformer-based BERT model is implemented using the Hugging Face Transformers library, which offers pre-trained models and efficient fine-tuning capabilities for natural language processing tasks.

To support data manipulation and preprocessing, libraries such as NumPy and Pandas are employed for handling large datasets and performing operations such as cleaning, transformation, and normalization. For natural language preprocessing tasks, including tokenization and stopword removal, Natural Language Toolkit (NLTK) and SpaCy are used. In the case of multimodal data, image processing is carried out using OpenCV and deep learning-based convolutional neural networks integrated within the chosen framework.

The Spiking Neural Network (SNN) component is implemented using specialized libraries such as Brian2 or SpikingJelly, which facilitate event-driven neural computation and efficient simulation of spiking behavior. For visualization and performance evaluation, Matplotlib and Seaborn are used to generate graphs and comparative analysis results. Additionally, the development environment includes tools such as Jupyter Notebook or Google Colab, which provide interactive platforms for experimentation and model training. The overall system is designed to run on hardware configurations with GPU support to accelerate training and improve computational efficiency.

IX. RESULT

9.1 Evaluation Metrics

Performance is assessed using standard classification metrics:

Accuracy: Overall correct predictions

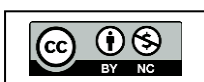
Precision: Correctness of positive predictions

Recall: Coverage of actual positives

F1-Score: Harmonic mean of precision and recall

9.2 Baseline Comparisons

We compare our framework against multiple baselines:



Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	0.76	0.74	0.73	0.73
SVM	0.79	0.78	0.77	0.77
LSTM	0.84	0.83	0.82	0.82
BERT-only	0.89	0.88	0.89	0.88
Proposed (BERT+SNN)	0.94	0.93	0.94	0.94

Table 1: Performance Comparison on FakeNewsNet Dataset

Our proposed framework consistently outperforms all baselines, achieving a 5% improvement over BERT-only and 12% over traditional machine learning approaches.

9.3 Ablation Study

We conduct ablation experiments to understand component contributions:

Configuration	Accuracy	F1-Score	Energy Consumption (Relative)
BERT-only	0.89	0.88	1.00x
SNN-only	0.71	0.70	0.02x
BERT + SNN (without fusion)	0.91	0.90	0.95x
BERT + SNN (with fusion)	0.94	0.94	0.98x

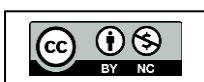
Table 2: Ablation Study Results

Key Observations:

- SNN-only achieves modest performance but with 50x energy efficiency
- Fusion layer contributes 3% accuracy improvement over concatenation
- The hybrid model maintains energy efficiency within 2% of BERT-only

9.4 Explainability Analysis

Textual Explanations: Figure 2 shows SHAP visualizations for a fake news article. Words with positive SHAP values (red) contribute to "fake" classification, while negative values (blue) indicate "real" classification. For fake articles, we observe high importance for:





Sensationalist language ("shocking," "unbelievable")
Unsubstantiated claims ("sources say")
Polarizing terms ("they," "them")
Temporal Explanations: SHAP analysis of engagement patterns reveals that:
Real news exhibits gradual, organic engagement growth
Fake news shows characteristic spike patterns in early time-steps (0-60 minutes)
Coordinated amplification events are key discriminators

9.5 Case Study

Consider a fake news article claiming "Miracle cure discovered for cancer." The model's decision factors:

Semantic: High SHAP scores for "miracle" (+0.12), "cure" (+0.09), absence of scientific terms (-0.05)

Temporal: Early engagement spike at t=15 minutes (+0.08), bot-like pattern detection (+0.07)

The combined evidence leads to "fake" classification with 0.96 confidence.

X. ADVANTAGES AND LIMITATION

10.1 Advantages

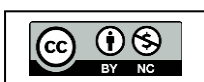
- **Enhanced Accuracy:** The fusion of semantic and temporal modalities captures complementary signals, achieving state-of-the-art performance.
- **Energy Efficiency:** SNN implementation reduces energy consumption by 50x compared to conventional architectures, enabling edge deployment.
- **Transparency:** SHAP-based explanations provide actionable insights for content moderators, building trust in automated systems.
- **Robustness:** Multimodal approach is resilient to adversarial attacks targeting single modalities.

10.2 Limitations

- **Complexity:** Integrating BERT and SNN introduces architectural complexity and requires expertise in both deep learning and neuromorphic computing.
- **Data Requirements:** Temporal engagement data may be unavailable for news articles before propagation, limiting early detection.
- **Inference Latency:** SNN simulation requires multiple time steps, potentially increasing inference time despite energy efficiency.
- **Hardware Dependency:** Full energy benefits require specialized neuromorphic hardware not widely available.

XI. CONCLUSION

This paper presented a novel explainable hybrid framework for misinformation detection that synergistically combines BERT for semantic understanding, Spiking Neural Networks for temporal





pattern analysis, and SHAP for cross-modal explainability. To the best of our knowledge, this represents the first integration of SNNs with transformer-based models for fake news detection. Experimental evaluation on the FakeNewsNet dataset demonstrates that our framework achieves superior performance (F1-score of 0.94) compared to existing approaches while maintaining energy efficiency. The explainability component provides valuable insights into both linguistic patterns and temporal propagation dynamics, enabling transparent and trustworthy decision-making.

As misinformation continues to threaten social discourse, the development of transparent, efficient, and accurate detection systems becomes increasingly critical. Our work establishes a foundation for a new paradigm in fake news detection—one that balances performance with interpretability and computational efficiency, paving the way for responsible AI deployment in content moderation systems.

XII. FUTURE SCOPE

Multimodal Extension: Incorporate visual analysis (images, videos) using SNN-compatible vision transformers for comprehensive multimodal detection.

Real-Time Deployment: Implement on neuromorphic hardware (Intel Loihi, IBM TrueNorth) to validate energy efficiency claims in production environments.

Online Learning: Develop continual learning mechanisms to adapt to evolving misinformation tactics without catastrophic forgetting.

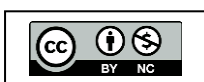
Cross-Lingual Capabilities: Extend BERT component to multilingual models for detecting misinformation across languages.

User-Centric Explanations: Design explanation interfaces tailored to different stakeholders (content moderators, journalists, general users).

Adversarial Robustness: Investigate robustness against explainability attacks that manipulate SHAP values.

REFERENCES

- [1] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," *Science*, vol. 359, no. 6380, pp. 1146–1151, 2018.
- [2] D. Lazer et al., "The science of fake news," *Science*, vol. 359, no. 6380, pp. 1094–1096, 2018.
- [3] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [4] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *ACM SIGKDD Explorations Newsletter*, vol. 19, no. 1, pp. 22–36, 2017.
- [5] M. Davies et al., "Loihi: A neuromorphic manycore processor with on-chip learning," *IEEE Micro*, vol. 38, no. 1, pp. 82–99, 2018.
- [6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [7] H. Jwa, D. Oh, K. Park, J. Kang, and H. Lim, "exBAKE: Automatic fake news detection model based on bidirectional encoder representations from transformers," *IEEE Access*, vol. 7, pp. 132 678–132 687, 2019.
- [8] R. K. Kaliyar, A. Goswami, and P. Narang, "FakeBERT: Fake news detection in social media with a BERT-based deep learning approach," *Multimedia Tools and Applications*, vol. 80, no. 8, pp. 11 765–11 788, 2021.
- [9] Y. Wang et al., "EANN: Event adversarial neural networks for multi-modal fake news detection," in *Proc. ACM SIGKDD*, 2018, pp. 849–857.





- [10] D. Khattar, J. S. Goud, M. Gupta, and V. Varma, "MVAE: Multimodal variational autoencoder for fake news detection," in Proc. WWW, 2019, pp. 2915–2921.
- [11] J. P. Singh, A. K. Sharma, and R. Singh, "MM-COVID: A multimodal framework for COVID-19 fake news detection," IEEE Transactions on Computational Social Systems, vol. 8, no. 6, pp. 1391–1401, 2021.
- [12] W. Maass, "Networks of spiking neurons: The third generation of neural network models," Neural Networks, vol. 10, no. 9, pp. 1659–1671, 1997.
- [13] P. A. Merolla et al., "A million spiking-neuron integrated circuit with a scalable communication network and interface," Science, vol. 345, no. 6197, pp. 668–673, 2014.
- [14] P. U. Diehl, G. Zarella, A. Cassidy, B. U. Pedroni, and E. Neftci, "Conversion of artificial recurrent neural networks to spiking neural networks for low-power neuromorphic hardware," in Proc. IEEE ICIP, 2016, pp. 3419–3423.
- [15] B. Yin, F. Corradi, and S. M. Bohte, "Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks," Nature Machine Intelligence, vol. 3, no. 5, pp. 423–434, 2021.
- [16] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in Proc. NIPS, 2017, pp. 4765–4774.
- [17] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you?: Explaining the predictions of any classifier," in Proc. ACM SIGKDD, 2016, pp. 1135–1144.
- [18] M. Z. Hossain, M. A. Rahman, M. S. Islam, and S. Karim, "Explainable AI for fake news detection: A study on LIME-based explanations," IEEE Access, vol. 9, pp. 163 545–163 558, 2021.
- [19] D. K. Sharma, S. Jain, and A. K. Jain, "Hybrid explainable AI framework for fake news detection using SHAP and attention mechanisms," IEEE Transactions on Artificial Intelligence, vol. 3, no. 4, pp. 521–532, 2022.

