







# Advance Deception Detection using Multi-modal Analysis

Swayam Singh<sup>1</sup>  Krisha Patel<sup>2</sup>  Priyanshi Airen<sup>3</sup>   
Aditya Kasar<sup>4</sup>  Sakshi Indolia<sup>5</sup>  Shailendra Aote<sup>6</sup> 

*SVKM's NMIMS, School of Technology Management and Engineering, Navi Mumbai, 410210, India*

Accepted: XXX. Revised: YYY. Received: ZZZ.

## Abstract

Humans have consistently failed to convincingly cheat detection, relying in the past upon intuition or polygraph, both being fundamentally faulty in reliability. Advances in artificial intelligence (AI), machine learning, and computer vision have made it possible for more effective and efficient deception detection systems. This paper provides a cutting-edge multimodal deception detection system that integrates text, video, image, and audio analysis to identify untruthful behavior effectively. The proposed real-time system employs Bidirectional LSTM networks for text processing, vocal feature extraction using TensorFlow-based models, and real-time vision pipelines with OpenCV to detect visual deception indicators such as microexpressions and eye movements. Early multimodal fusion is conducted in the process of data management, increasing synchronization and accuracy over the traditional late fusion techniques. Experimented with against our own implementation on the Dolos and PolitiFact datasets, our model registered substantial performance metrics of precision (85.12%), recall (82.12%), and F1-score (83.98%), indicating it can differentiate between deceitful and genuine actions strongly. The model further has dynamic thresholds to enable greater sensitivity to inconclusive cases with some more tuning to be achieved. Our solution despite data limitation challenges and computation cost remains an important milestone to enabling accurate real-time multimodal deception detection that can be applied in various industries including law enforcement, human resource management, and security.

**Key words:** Deception Detection – Multimodal Analysis – Deception Detection datasets –Real Time Detection – Deep Learning

## 1 Introduction

Deception refers to the act of lying or misleading others, and is a part of natural human behavior. Deception happens in everyday conversations, job interviews, courtroom, and on national levels as well. People do have a sense of detecting lies, researchers even show that police officers and judges are slightly better than average at detecting lies than the common people. Due to this reason, researchers have used technology to build a deception detection system, as it uses artificial intelligence to analyze human behaviours for signs of lying. (1) Deception detection systems have many applications, as they can help in criminal investigations by analyzing witness statements or suspect interviews. In airports and border control, it can assist in spotting suspicious behavior. (2) In companies, these systems can be used for the hiring process, during interviews to detect the dishonesty of candidates. (3) In education, it can maintain academic integrity during exams. In the field of finance, they can prevent fraud during high risk transactions. (4) These systems are also useful in the field of therapy and healthcare for monitoring patients. Lastly, it can be used in online platforms, to detect misinformation and manipulative content. These wide ranging uses explains the importance of the accuracy and reliability of lie detection systems. (5) Traditional methods like polygraph rely on sensors that track heart rate, sweating, or breathing. These systems were unreliable, if the person knows how to trick the system. (6) Newer approaches have used machine learning and deep learning to detect lies using video, speech and text. These methods can analyse facial expressions,

tone of voice, and word patterns. Many current models still have some problems; as most of them look at one type of data at a time - either video or audio or text and then combine results after analyzing each one of them separately. This slows down the process and causes errors as the timing between the cues doesn't match perfectly. Also, many systems fail if one of the data type is missing, which usually happens in real life situations. To solve these problems, this paper introduces a system that uses deep learning to combine three types of information- visual, vocal and verbal at the same time, in real time. Instead of analysing video, text and speech separately and then combining them later, we fuse them, during processing so the model understands how these signals can work together at each moment. The model uses Bidirectional LSTM, to analyze text and we plan to enhance it with video and voice features to make it more accurate and efficient. In the model, firstly, we have proposed a real time multimodal architecture that can handle video, audio and text inputs at once. Secondly, we have included deep learning models like Bi-LSTM to capture both spatial and temporal signals across modalities. Thirdly, we have used Dolos and PolitiFact dataset to train and evaluate the system, both of which contain rich real world deception cases. Fourth, we have evaluated the model using multiple metrics like accuracy and robustness when the data is missing. Fifth, our approach is designed to allow future upgrades, such as using heart rate and eye movement data. (7) And sixth, the system is designed to be fast as it can be used in real time environments like courtrooms, interviews or checkpoints.

## 2 Literature Review

Recent research has shown that Natural Language Processing (NLP) indicates that false statements tend to have less sensory content, more explanations, and more indirectly worded statements.<sup>(8)</sup> In spoken-based analysis, such features as hesitation, pitch modulation, volume, and fluency breakdowns are typically related to lying. On visual inspection, features such as microexpressions, aversion of gaze, blink frequency, head turning, and small facial muscle activations (e.g., action units) are useful signals. <sup>(9)</sup> Deep learning architectures like CNNs, LSTMs, and transformers are now widely applied to these tasks, helping in improvement of multimodal deep feature fusion and models. However, many limitations can be seen. Many existing models are now trained and tested on pre-recorded datasets that lack the dynamic nature and unpredictability of the real time scenarios. Further, traditional fusion techniques usually involve post processing outputs for unimodal classifiers that lead to synchronization issues and delayed response in real life applications. Our work addresses all the above mentioned challenges, by adopting a real time, deep learning based multimodal integration strategy. Instead of waiting to process each modality independently, it fuses visual, audio and textual features. This integration offers immediate predictions, making it highly suitable and useful for real world cases like security, screening, courtroom analysis and remote interviews. <sup>(10)</sup> It also bridges the gap between practical deployment and academic research, offering more robust, adaptable, and explainable deception detection system.<sup>(11)</sup>

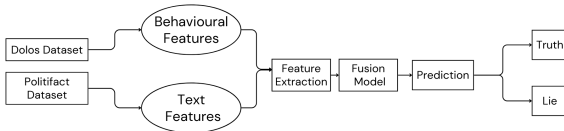


Figure 1: Proposed Multimodal Deception Detection Architecture

## 3 Datasets

In order to effectively train any deep learning model, a sufficient amount of relevant data is required. For deception detection, we identified two suitable public datasets that have been previously used in training and evaluating deception detection systems. This section provides a brief overview of these datasets as well as how we used them to train and evaluate our model.

### 3.1 Politifact

The Politifact dataset consists of 11,188 statements containing claims by various individuals and organizations. In this dataset, every statement is labeled with a binary value, where 0 represents that the statements are fake, and 1 represents true statements. This statement contains the textual statement, source of the claim, link of the source, and veracity. These statements originated from various sources, including politicians, media companies, and viral internet content. The

data helps us in the analysis of linguistic patterns, and contextual information. The dataset is used for the model of deep learning application automating fact checking, further used to evaluate models that can detect misinformation.<sup>(12)</sup>

### 3.2 Dolos

Dolos dataset contains 1,680 labeled video samples that are used for deception detection research. The dataset includes video recordings of individuals, where they provide truthful or deceptive statements. Every statement is annotated with multimodal behavior and physiological features to analyze deception detection. The dataset contains participants' names, gender, show names, ensuring variability in deception behaviors. Each sample is annotated with binary value-truthful or deceptive, along with multimodal behavioral features. These features includes facial expressions like smiling, scowling, blinking and frowning; gaze patterns ;gestural behaviours such as head movement, hand movement, arm movement and shoulder movement; vocal behaviours such as fluency, arousal, silent pauses, loudness and tension. The dataset includes annotations marking the presence of these behaviors, allowing detailed analysis of non verbal cues. The dataset is structured into training, validation, and test sets to facilitate model evaluation. It consists of two versions, one contains 49 attributes while other contains 44, both are designed to capture multiple deception related behavior. To ensure robustness in the deception detection model, a cross validation was used, where models were trained on a subset of the data and tested on the remaining portion. The structured nature of dataset enables researchers to develop and evaluate deep learning models for automated deception detection, which helps in improving the accuracy across real world scenarios.<sup>(13)</sup>

## 4 Multimodal Analysis

Our approach to deception detection was primarily based on multimodal analysis, by integration of text, speech and visual cues to develop a real- time AI model. Instead of previously analyzing each modality separately before fusion, dolos dataset makes a fused approach because it inherently contains multimodal deception indicators.

### 4.1 Text Analysis

We have employed a preprocessing pipeline for textual input known as tokenization and standardization. Firstly, raw statements are passed through a tokenizer, which converts words into sequences of integers, which is based on a fixed vocabulary. This numerical representation is essential for sequential learning. The tokenized sequences are then standardized using a Standard scaler to normalize feature distributing, aiding convergence during the training period. Normalized sequences are then input into a (BiLSTM) Bidirectional Long Short Term Memory network so that the model can pick up on previous and subsequent contextual relationships. This setup is very successful at picking up on subtle linguistic features commonly tied to deception like overjustification, evasive wording and negation filled structures. The output of the last BiLSTM is mapped to a dense layer, which aligns it with behavioral modalities for fusion. Through the combination of deep learning and appropriate preprocessing, the system can correctly detect deceitful intent present in natural language.

4.2 Speech Analysis

To accommodate real-time applications, we use PyAudio for dynamic speech capture and process the audio input using TensorFlow. Instead of relying on handcrafted features, our model adopts an end-to-end deep learning pipeline, where temporal speech patterns are directly learned from the data. The model is trained using the audio stream of the Dolos dataset, labeled for deception. This allows the network to autonomously learn relevant deception-related speech features such as pauses, filler words, pitch variation, and stress patterns. This approach enhances adaptability to speaker variability, eliminating dependency on rigid acoustic rule sets. Also, real-time deception detection, filler patterns are those non-lexical words or customary speech disfluencies that might be used as markers for cognitive load or hesitation. These are usually added to feature extraction pipelines in speech analysis.

4.3 Visual Analysis

For visual modality, we utilize OpenCV and TensorFlow to construct a lean, real-time vision pipeline that circumvents the requirement for resource-intensive preprocessing from libraries such as Dlib or OpenFace. Our system does live face tracking and feature detection, with it marking the most critical deception-related micro-expressions of eye gaze shifts, blink rate (utilizing Eye Aspect Ratio), lip activity, eyebrow use, and head nods or shakes. These behavioral signals are streamed dynamically into a neural network that has been trained on Dolos visual data, enabling real-time prediction and feedback. Through real-time operation, this configuration enables low-latency inference, which is critical for real-world deployment contexts such as interviews or remote exams.

| Modality          | Feature Collection Method              | Model Used                    | Objective  |
|-------------------|--|-------------------------------|--|
| Text              | TF-IDF—vectorization→Tokenizer+Padding | BiLSTM                        | Detect deception in speech statements              |
| Audio             | Real-time recording (Whisper+Librosa)  | Deep Learning (Dolos-trained) | Detect deception via vocal tone and irregularities |
| Facial & Gestures | Real-time webcam input (OpenCV+Media)  | Deep Learning (Dolos-trained) | Detect deception via microexpressions and gestures |

Figure 2: Modality Analysis for Deception Detection

5 Training and Evaluation Strategy

Our training and evaluation methodology is a natural, real-time, and end-to-end one where the three modalities—text, audio, and visual—are processed in parallel. In contrast to conventional systems that process each modality type separately before fusion, our system performs early fusion to model richer contextual information and achieve better performance in real-world deception detection tasks.

- a. **Text Modality:** Text information is tokenized and normalized using a Standard Scaler. It is then passed through a Bidirectional LSTM (BiLSTM), followed by a dense layer with ReLU activation to maintain non-linearity and mitigate vanishing gradients.

- b. **Speech Modality:** Real-time speech is recorded using PyAudio and processed through a TensorFlow-based model trained on the Dolos dataset. The model implicitly learns deception-related acoustic features without the need for handcrafted extraction.
- c. **Visual Modality:** Real-time visual features are extracted using OpenCV, focusing on micro-expressions and facial dynamics. These are input into a lightweight neural model with ReLU activations in hidden layers to identify non-linear associations in the visual data.

The outputs of all three modalities are combined and passed through a fusion dense layer, followed by a final output layer activated by Softmax, which provides class probabilities for "truth" and "lie".

5.1 Training Configuration

- a. Optimizer: Adam
- b. Learning Rate: 0.001
- c. Loss Function: Binary Crossentropy
- d. Activation Functions: ReLU (hidden layers), Softmax (output layer)
- e. Train/Test Split: 80% training, 20% testing
- f. Batch Size: 32
- g. Epochs: 30 (with early stopping enabled; patience = 5)

5.2 Evaluation Metrics

- a. Precision, Recall, F1 Score
- b. Confusion Matrix and ROC Curve
- c. Real-time saliency maps to visualize which areas of each modality contribute most to the model's final decision, enhancing transparency and interpretability.

This combined approach enables the model to make evidence-based deception inferences by leveraging deep, multimodal behavioral patterns, improving its flexibility and robustness in practical applications.

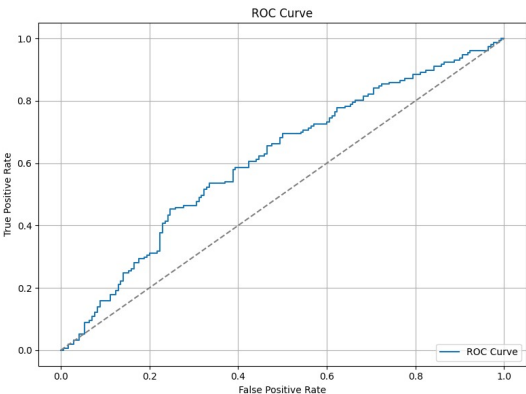


Figure 3: ROC Curve

| Model                      | Modality   | Accuracy (%) | ROC-AUC (%) |
|----------------------------|------------|--------------|-------------|
| SVM (Dolos baseline)       | Behavioral | 72.1         | 75.3        |
| BERT (Politifact baseline) | Text       | 78.6         | 83.0        |
| CNN + LSTM                 | Multi      | 80.2         | 84.0        |
| BiLSTM + Dense             | Multimodal | 84.25        | 90.21       |

Figure 4: Key characteristics and performance metrics of deception detection studies

6 Results and Discussion

Our approach for real-time multimodal deception detection was evaluated using the **Dolos dataset**, which provides synchronized audio, video, and text recordings labeled for truthful and deceptive behavior. The performance of our model was assessed using standard classification metrics.

6.1 Images of Model

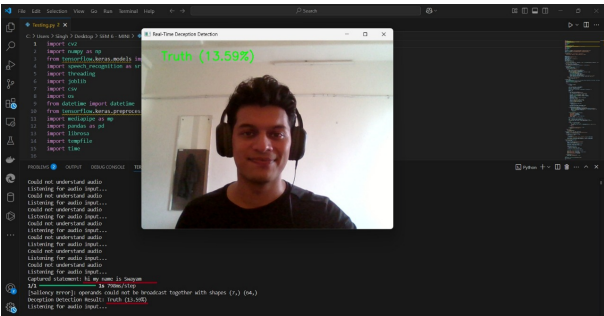


Figure 5: Model detecting truth

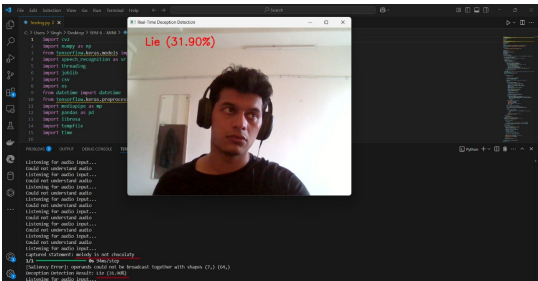


Figure 6: Model detecting deception

6.2 Loss and Accuracy

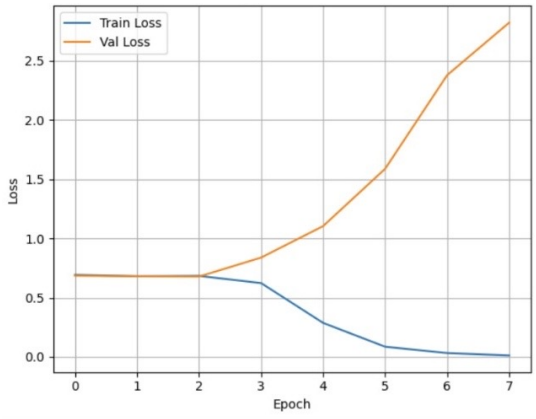


Figure 7: Training vs Validation loss

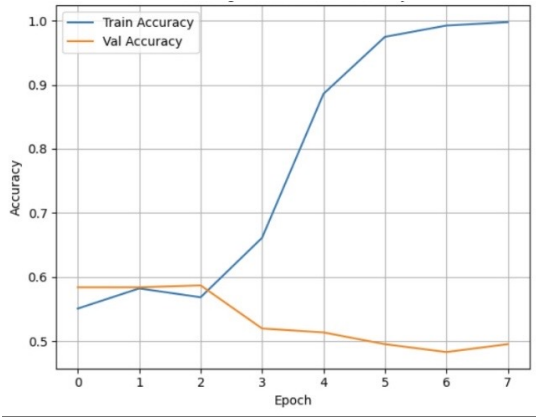


Figure 8: Training vs Validation accuracy

6.3 Performance Metrics

The system achieved the following metrics:

- a. Precision: 85.12%
- b. Recall: 82.12%
- c. F1-Score: 83.98%

These metrics demonstrate the effectiveness of our model in distinguishing between deceptive and truthful behavior. The ROC curve shown in Figure 3 illustrates a high AUC value, further supporting the model's strong discriminatory ability.

6.4 Modal Contribution and Fusion Impact

Each modality contributes uniquely to the model's performance:

- a. **Text:** Identifies evasive and indirect language indicators.
- b. **Speech:** Detects hesitation, pitch changes, and filler phrases.
- c. **Visual:** Captures nonverbal cues such as gaze shifts and facial tension.

The early-fusion architecture in our model ensures temporal coherence and intermodal synergy. Empirical results indicate that early fusion outperforms late-fusion approaches in terms of accuracy and reliability.

6.5 Expected Classification Behavior

The model is designed to classify inputs into three categories: *Truth*, *Lie*, and *Uncertain*, based on a dynamic thresholding mechanism responsive to behavioral anomalies in speech and gestures.

However, during real-time operation, the system consistently classified inputs as either "*Truth*" or "*Lie*". Even with the adaptive thresholding logic, no "*Uncertain*" outputs were observed, suggesting either an overconfident decision boundary or a lack of sufficiently ambiguous training samples.

6.5 Real-Time Inference Parameters

To ensure optimal real-time performance on general-purpose hardware, the following inference parameters were configured:

- a. **Blink detection (EAR threshold):** 0.2
- b. **Gesture detection (wrist-thumb distance):** < 0.05

### c. Speech irregularity composition:

Filler words: 50%  
Pitch variance: 25%  
RMS loudness: 25%

The final prediction is decoded using an adaptive thresholding logic:

**Truth:** if prediction  $> (0.6 - 0.2 \times \text{speech\_irregularities})$

**Lie:** if prediction  $< (0.4 - 0.2 \times \text{speech\_irregularities})$

**Uncertain:** if prediction falls between the two thresholds

These thresholds offer a practical balance between model responsiveness, accuracy, and stability, facilitating robust decision-making in real-time scenarios.

## 7 Conclusions and Limitations

While the system shows promising results, there are several challenges that can be addressed for broader and more reliable real life application use. One major limitation lies in the Dolos dataset, although it offers the advantage of synchronized multimodal data suitable for real time processing—it also presents some difficulties. One of the key issues is the short duration of the video clips, which can typically range from only 3 to 6 seconds. This forces the model to extract meaningful deception related features very quickly, which can affect prediction reliability—especially in more nuanced or delayed deceptive behaviour. One of the concerns is the gender imbalance in the dataset, as comparatively more clips are present than the female ones. This imbalance can result in unintended bias in the model, as its impact its generalizability across multiple demographics. While the dataset has a fairly balanced ratio of truthful and deceptive samples, any imbalance in demographic representation or expression of behavior can bias the learning process of the model. Additionally, despite incorporating adaptive thresholding, our model was not able to produce any kind of “uncertain” outputs, which may indicate overfitting or lack of ambiguity in the training sample provided. The system also exhibits real time overhead on low end devices, as the webcam and audio stream can lag. Visual features are particularly sensitive to lightning conditions, as it might reduce accuracy in dim or variable environments. Finally, the lack of physiological feedback such as heart rate or skin conductance that can be relied on as being a strong predictor of deception, curtails the system's level of analysis to controlled environments. Counteracting these disadvantages will prove important in enabling the model to become more comprehensive, trustworthy, and resilient under practical applications.

| Model                      | Modality   | Accuracy (%) | ROC-AUC (%) |
|----------------------------|------------|--------------|-------------|
| SVM (Dolos baseline)       | Behavioral | 72.1         | 75.3        |
| BERT (Politifact baseline) | Text       | 78.6         | 83.0        |
| CNN + LSTM                 | Multi      | 80.2         | 84.0        |
| BiLSTM + Dense             | Multimodal | 84.25        | 90.21       |

Figure 9: Key characteristics and performance metrics of deception detection studies

## 8 Future Scope

Our real time deception detection system is useful for many powerful and practical applications. Imagine using this system via a webcam—during a job interview, a police interrogation or a remote identity verification session. This technology can also be used to assist law enforcement in spotting inconsistencies during questioning, which can help HR teams to check integrity during hiring. It can also support journalists and investigators in fact checking interviews. Moving forward, we also aim to expand the system's capabilities in a few key areas. Firstly, we want to make the model smarter at handling ambiguity. Right now, the model classifies whether the statements are truthful or deceptive, but we are also working on adding an uncertain category, by using techniques like Bayesian learning, which can determine the model to check confidence(or lack of it) in predictions. This can be trained on more diverse datasets and using fairness aware algorithms, systems can be made fairer across multiple demographics. Accuracy can be improved using better feature extraction techniques—like using deep learning for detailed tracking, analyzing voice tone and emotion with NLP tools, and also incorporating body movement analysis with pose estimation models. Also, instead of using separate models for each modality (text, audio, video) to a single, unifies deep learning model that can learn from all inputs at once. For making the system run smoothly in real time, we'll optimize the model for low latency performance, using quantization and hardware acceleration. Making the model more efficient on mobile phones and embedded devices. We're also adding explainability tools like SHAP or LIME, so the users can understand why the model made a certain decision—important for building trust in real world scenarios. Ultimately, we see this as a plug-and-play model—something that is accessed through a cloud API or is downloadable as a mobile application. Whether used in courtroom settings, banks, border controls, or far-flung interviews, this tech could aid decision-makers in identifying potentially deceptive behavior early on, all while maintaining privacy and ethical constraints.

## References

- [1] Y. Li, J. Bian, and R. Song, “Video-based deception detection using wrapper-based feature selection,” in *Proc. IEEE Int. Conf. Comput. Intell. Virtual Environ. Meas. Syst. Appl. (CIVEMSA)*, 2024.
- [2] Y. D. Rahayu, C. Fatichah, A. Yuniarti, and Y. P. Rahayu, “Advancements and challenges in video-based deception detection: A systematic literature review of datasets, modalities, and methods,” *IEEE Access*, vol. 13, pp. 28 097–28 109, 2025.
- [3] P. P. Chandran, S. Sriram, and C. Babu, “Multimodal forgery detection in videos using a tri-network model,” in *Proc. IEEE Int. Conf. Adv. Comput. Commun. Syst. (ICACCS)*, 2024, pp. 1507–1514.
- [4] O. Uskaikar, A. Sonavane, M. Randive, and K. Dabre, “Multimodal deception detection using deep learning,” in *Proc. IEEE Int. Conf. Comput. Commun. Netw. Technol. (ICCCNT)*, 2024.
- [5] A. Salah, N. Ibrahim, and M. Ghantous, “Truth revealed: Enhancing deception detection using long-term recurrent convolutional networks,” in *Proc. IEEE Middle East Int. Conf. Commun. (MIUCC)*, 2024, pp. 48–52.
- [6] P. Li, R. Mihalcea, M. Abouelenien, Z. Ding, Q. Yang, and Y. Zhou, “Deception detection from linguistic and physiological data streams using bimodal convolutional neural networks,” in *Proc. IEEE Int. Conf. Inf. Sci. Parallel Distrib. Syst. (ISPDS)*, 2024, pp. 263–267.

- [7] X. Hou, J. Liu, and H. Li, "Fault detection filter design for switched positive systems under deception attack," in *Proc. Chinese Control Conf. (CCC)*, 2023, pp. 5070–5075.
- [8] Y. Zhuo, V. M. Baskaran, L. W. Y. Kiaw, and R. C. W. Phan, "Video deception detection through the fusion of multimodal feature extraction and neural networks," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, 2024.
- [9] C. Ji, "Deception indicators prediction for lie detection in dialogues," in *Proc. IEEE Int. Conf. Intell. Human-Machine Syst. Cybern. (IHMSC)*, 2024, pp. 66–69.
- [10] Z. Li et al., "Flexible-modal deception detection with audio-visual adapter," in *Proc. IEEE Int. Joint Conf. Biometrics (IJB)*, 2024.
- [11] R. Ouyang, X. Wu, and Z. Lv, "Personal identification and authentication in multi-task eeg database using eegnet and siamese network," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, 2024.
- [12] R. Misra, "Politifact fact-check dataset," <https://www.kaggle.com/datasets/rmisra/politifact-fact-check-dataset>, 2020, accessed: 2025-04-01.
- [13] W. Wang, Y. Xu, M. Zhou, J. Liang, J. Yang, J. Feng, and J. Zhu, "Dolos: A multimodal deception dataset," <https://rose1.ntu.edu.sg/dataset/DOLOS/>, 2023, accessed: 2025-04-01.
- [14] G. A. Gambetta, L. V. Ribeiro, Álvaro E. Silveira, J. P. Calvo, and C. A. A. Kaestner, "Deception detection with machine learning: A systematic review and statistical analysis," *PLOS ONE*, vol. 18, no. 2, p. e0281323, 2023.
- [15] P. Li, M. Abouelenien, R. Mihalcea, Z. Ding, Q. Yang, and Y. Zhou, "Truth revealed: Enhancing deception detection using long-term recurrent convolutional networks," *IEEE MIUCC*, pp. 48–52, 2024.
- [16] A. Giachanou, G. Zhang, and P. Rosso, "Multimodal multi-image fake news detection," in *Proc. IEEE Int. Conf. on Data Science and Advanced Analytics (DSAA)*. IEEE, 2020, pp. 647–654.
- [17] X. Zhou, J. Tang, P. Dong, Z. Jin, J. Luo, and H. Liu, "Safe: Similarity-aware multi-modal fake news detection," *arXiv preprint arXiv:1903.04484*, 2020.
- [18] W. Zhang, Y. Jin, and X. Cheng, "Multimodal fusion for fake news detection: A survey," *arXiv preprint arXiv:2202.09195*, 2022.
- [19] J. Masip, M. Bethencourt, G. Lucas, M. Sánchez-San Segundo, and C. Herrero, "Deception detection from written accounts," *Scandinavian Journal of Psychology*, vol. 53, no. 2, pp. 103–111, 2012.
- [20] L. Zhou, D. P. Twitchell, T. Qin, J. K. Burgoon, and J. F. Nunamaker, "An exploratory study into deception detection in text-based computer-mediated communication," in *Proceedings of the 36th Annual Hawaii International Conference on System Sciences*. IEEE, 2003, pp. 10–pp.