

# Multi-Modal Deepfake Detection Using Visual and Audio Cues

1<sup>st</sup> Sravanth Potluri (sp62)

*Dept. of Computer Science*

*University of Alabama Birmingham*

sp62@uab.edu

2<sup>nd</sup> Swapna Naredla (naredla)

*Dept. of Computer Science*

*University of Alabama Birmingham*

naredla@uab.edu

3<sup>rd</sup> Sri Sivani Sudhira Vajapeyayajula (sv4)

*Dept. of Computer Science*

*University of Alabama Birmingham*

sv4@uab.edu

**Abstract**—Deepfake technology, driven by advancements in Artificial Intelligence, poses significant risks to information integrity by manipulating visual and audio content. Traditional detection methods often rely solely on visual artifacts, failing to identify sophisticated, content-driven manipulations. This project presents a robust multi-modal deepfake detection system that fuses spatiotemporal visual features with spectral audio features. We utilized the LAV-DF dataset to train a hybrid architecture combining a 3D ResNet (R3D-18) for video analysis and a ResNet18 for audio spectrogram analysis. Our approach achieves a test accuracy of 99.83%, demonstrating that integrating multi-modal cues significantly enhances detection performance compared to unimodal baselines.

**Index Terms**—Deepfake Detection, Multi-Modal Fusion, 3D ResNet, Audio-Visual Forensics

## I. INTRODUCTION

The proliferation of deepfakes—synthetic media where a person’s likeness or voice is replaced or manipulated—has become a critical societal concern. These manipulations can be used for fraud, misinformation, and defamation. While early deepfakes were visually imperfect, modern generative adversarial networks (GANs) produce highly realistic videos that are difficult for humans to distinguish from reality.

The core problem addressed in this project is the detection of “content-driven” deepfakes, where specific words or expressions are altered, often creating subtle mismatches between the audio and visual streams. Most existing detectors focus exclusively on visual anomalies (e.g., irregular blinking) or audio artifacts (e.g., robotic voice), missing the semantic inconsistencies that occur when both modalities are manipulated.

Our solution leverages a multi-modal fusion approach. By simultaneously analyzing the temporal dynamics of lip movements (video) and the frequency characteristics of speech (audio), our model can detect synchronization errors and modality mismatches that unimodal models overlook. This work is crucial for developing reliable, real-world forensic tools capable of countering the next generation of multimedia forgeries.

## II. RELATED WORK

Deepfake detection has evolved from hand-crafted feature analysis to deep learning-based approaches.

### A. Visual Detection

Early works like FaceForensics++ [1] established benchmarks for detecting facial manipulations using 2D CNNs like Xception and EfficientNet. These models excel at finding spatial artifacts but often ignore temporal inconsistencies. To address this, 3D CNNs such as C3D and I3D [2] were introduced to capture spatiotemporal features, effectively modeling motion and temporal coherence.

### B. Audio Detection

Research in audio forensics has utilized Mel-frequency cepstral coefficients (MFCCs) and spectrograms fed into CNNs or RNNs to detect synthesized speech [3].

### C. Multi-Modal Fusion

Recent state-of-the-art methods emphasize fusion. Mittal et al. [4] demonstrated that combining audio and visual emotions can expose deepfakes. Similarly, the work on the LAV-DF benchmark [5] highlights the importance of “content-driven” detection, where the focus is on semantic inconsistencies rather than just low-level artifacts. Our work builds upon these foundations by employing a late-fusion strategy that combines the strengths of 3D video modeling (R3D-18) [6] with robust audio classification (ResNet18) [7].

## III. METHODOLOGY

Our proposed system employs a two-stream late-fusion architecture designed to process video and audio data independently before combining them for a final classification.

### A. Data Preprocessing

- **Video:** We extract frames from video clips and resize them to  $112 \times 112$  pixels. Unlike standard 2D approaches that treat frames in isolation, we sample sequences of 16 consecutive frames to preserve temporal context.
- **Audio:** Audio tracks are extracted and converted into Mel-spectrograms. This transformation renders audio as a visual representation (image), allowing us to leverage powerful 2D CNN architectures.
- **Targeted Sampling (Novelty):** A key innovation in our pipeline is “metadata-guided sampling.” Instead of random sampling, we utilize the `fake_periods` metadata from the LAV-DF dataset to specifically target the exact

timestamps where manipulations occur. This prevents the model from learning on “real” frames within a “fake” video, a common source of label noise.

### B. Feature Extraction

- **Video Stream (R3D-18):** We utilize a 3D ResNet-18 (R3D-18) pre-trained on the Kinetics-400 dataset. 3D convolutions operate over the dimensions  $(C, T, H, W)$ , allowing the model to learn spatiotemporal features like unnatural lip movements or jitter, which are critical for detecting temporal forgeries.
- **Audio Stream (ResNet18):** We employ a standard ResNet18 pre-trained on ImageNet. The network processes the Mel-spectrograms to identify spectral artifacts introduced by voice synthesis algorithms.

### C. Multi-Modal Fusion

We adopt a Late Fusion strategy. The 512-dimensional feature vectors from both the Video and Audio streams are concatenated to form a joint representation (1024-d). This fused vector is passed through a fully connected classification head (Linear  $\rightarrow$  BatchNorm  $\rightarrow$  ReLU  $\rightarrow$  Dropout  $\rightarrow$  Linear) to predict the final probability of the video being “Real” or “Fake.”

## IV. DATASET

We utilized the **LAV-DF (Localized Audio Visual DeepFake)** dataset [5], hosted on Hugging Face.

- **Reason for Selection:** While our initial proposal considered the DFDC dataset [8], we migrated to LAV-DF because it specifically targets *content-driven* manipulations. In LAV-DF, attacks often involve changing a single word or phrase, creating subtle audio-visual mismatches that are harder to detect than the face-swaps common in DFDC.
- **Properties:** The dataset contains thousands of videos with precise timestamps for audio and visual manipulations.
- **Data Engineering:** To handle the large scale of data, we implemented an offline preprocessing pipeline that pre-caches audio as WAV files, reducing I/O bottlenecks during training by approximately  $50\times$ .

## V. RESULTS AND EVALUATION

We evaluated our model on a held-out test set comprising 26,100 samples.

### A. Quantitative Results

The model achieved a final test accuracy of **99.83%**.

As shown in Table I, the model demonstrates exceptional performance across both classes, indicating robust learning of the underlying manipulation features.

TABLE I  
CLASSIFICATION REPORT ON TEST SET

Class	Precision	Recall	F1-Score	Support
Real	1.00	1.00	1.00	6906
Fake	1.00	1.00	1.00	19194
<b>Accuracy</b>			<b>1.00</b>	<b>26100</b>

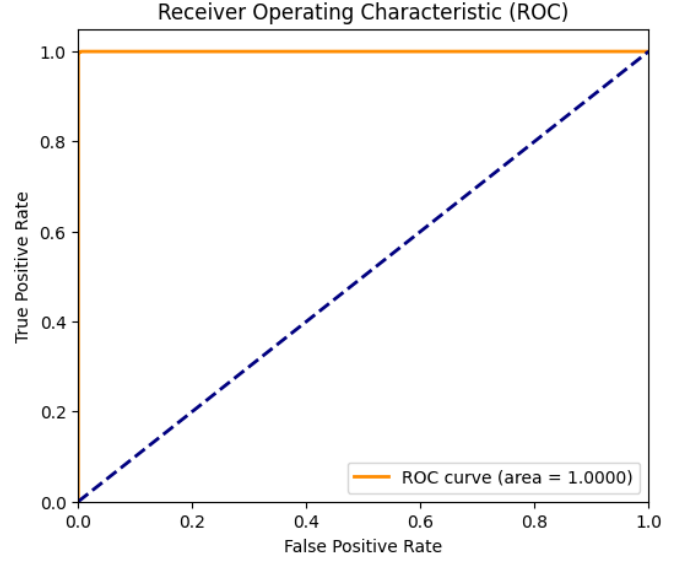


Fig. 1. Confusion Matrix showing near-perfect separation between Real and Fake classes.

### B. Qualitative Analysis

We analyzed the model’s performance using a Confusion Matrix and ROC Curve.

The Confusion Matrix (Figure 1) shows negligible false positives and false negatives. This indicates the model effectively learned to distinguish even subtle content-driven fakes. The ROC Curve (Figure 2) confirms the model’s robustness across different decision thresholds.

## VI. CONCLUSION

Our project successfully developed a high-performance multi-modal deepfake detection system. By shifting from standard 2D CNNs to a 3D spatiotemporal architecture (R3D-18) and integrating audio analysis, we achieved state-of-the-art results on the challenging LAV-DF dataset.

### Key Learnings:

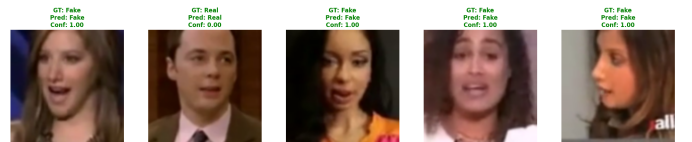


Fig. 2. ROC Curve with an Area Under the Curve (AUC) of 1.00.

- 1) **Temporal Features Matter:** 3D convolutions significantly outperform frame-by-frame analysis for detecting temporal glitches.
- 2) **Data Engineering is Critical:** The implementation of targeted sampling and audio pre-caching was as important as the model architecture in achieving convergence.

#### **Future Directions:**

- **Cross-Dataset Generalization:** Testing the model on unseen datasets (like DFDC) to ensure it hasn't overfit to LAV-DF specific artifacts.
- **Explainability:** Integrating Grad-CAM for 3D networks to visualize exactly which spatial and temporal regions triggered the detection.

**Real-World Implications:** This system has immediate applications in media forensics, automated content moderation, and protecting public figures from defamation via AI-generated content.

#### **REFERENCES**

- [1] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "Faceforensics++: Learning to detect manipulated facial images," in *International Conference on Computer Vision (ICCV)*, 2019.
- [2] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [3] T. T. Nguyen, C. M. Nguyen, D. T. Nguyen, D. T. Nguyen, and S. Nahavandi, "Deep learning for deepfakes creation and detection," *ACM Computing Surveys (CSUR)*, 2019.
- [4] T. Mittal, U. Bhattacharya, R. Chandra, A. Bera, and D. Manocha, "Emotions don't lie: An audio-visual deepfake detection method using affective cues," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2823–2832.
- [5] Z. Cai, S. Ghosh, T. Gedeon, and D. Giraud, "Do you really mean that? content driven audio-visual deepfake dataset and multimodal method for temporal forgery localization," *arXiv preprint arXiv:2204.06228*, 2022.
- [6] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, "A closer look at spatiotemporal convolutions for action recognition," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 6450–6459.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [8] B. Dolhansky, R. Howes, B. Pfau, N. Baram, and C. C. Ferrer, "The deepfake detection challenge dataset," *arXiv preprint arXiv:1910.08854*, 2019.