

| ISSN: 2347-8446 | www.ijarcst.org | editor@ijarcst.org |A Bimonthly, Peer Reviewed & Scholarly Journal

||Volume 6, Issue 2, March-April 2023||

DOI:10.15662/IJARCST.2023.0602001

Deepfake Detection using Multimodal Machine Learning Techniques

Anuj Dhar

Shree Ramchandra College of Engineering, Pune, India

ABSTRACT: Deepfake media—synthetic audiovisual content crafted using advanced AI models—pose serious threats to trust, security, and privacy. Traditional detection methods focusing solely on visual cues are increasingly circumvented by more sophisticated forgeries. Multimodal machine learning, which combines audio and video modalities, offers a more resilient detection framework. This review surveys pre-2022 research on multimodal deepfake detection, highlighting approaches that exploit both emotional inconsistency and temporal misalignment between modalities. For example, Mittal et al.'s affective-cue-based Siamese network achieves AUC of 96.6% on the DeepFake-TIMIT dataset and 84.4% on DFDC by modeling emotion discrepancies across audio and visual channels arXiv. Khalid et al. evaluate unimodal, ensemble, and multimodal detectors over the FakeAVCeleb dataset, observing that neither unimodal nor naive multimodal baselines consistently outperform ensemble methods arXiv. Additional methods include temporal feature prediction leveraging contrastive learning to detect misuse of synchronization between audio-visual sequences, achieving accuracy ~84.3% and AUC ~89.9% on FakeAVCeleb MDPI. Survey studies reinforce that deepfake detection must span modalities to counter increasingly seamless forgeries PeerJ. We propose a general workflow: data collection \rightarrow synchronized feature extraction \rightarrow modality-specific embedding \rightarrow cross-modal consistency modeling (e.g., contrastive, affective) → fusion and classification → evaluation. Advantages of multimodal frameworks include resilience to single-modality manipulation and modeling of inter-modal anomalies; disadvantages include the complexity of aligning asynchronous content, limited multimodal datasets, and heavier computation. We conclude that multimodal detection adds robustness but remains under-explored. Future directions include better dataset creation, fine-grained artifact modeling, self-supervised pretraining, and leveraging large pretrained multimodal models.

KEYWORDS: Deepfake Detection, Multimodal Learning, Audio–Visual Consistency, Affective Cues, Temporal Feature Prediction, DeepFake-TIMIT, FakeAVCeleb

I. INTRODUCTION

Deepfake technology enables the realistic synthesis and manipulation of video and audio, threatening authenticity in media and enabling misinformation, impersonation, and fraud. Traditional deepfake detection techniques often rely on visual artifacts (e.g., frequency anomalies, facial landmarks, blinking irregularities), but as generative models advance, such cues become more subtle or imperceptible.

In response, multimodal machine learning—integrating both **audio and visual** modalities—offers a compelling alternative by evaluating **inter-modal consistency**. Humans naturally rely on audiovisual synchronization for credibility assessment; thus, modeling discrepancies (e.g., mismatched emotion, timing, lip-sync) enhances detection robustness. Mittal et al.'s study introduced an affective-cue-based Siamese network that measures emotion alignment between modalities, delivering high detection performance using DeepFake-TIMIT and DFDC datasets arXiv.

Another challenge is the scarcity of **comprehensive multimodal datasets**. The FakeAVCeleb dataset addresses this by offering synchronized deepfake audio—video samples. Khalid et al. benchmarked unimodal, ensemble, and multimodal detectors on FakeAVCeleb, finding that naïve multimodal approaches underperform compared to ensemble models arXiv.

Finally, temporal inconsistencies offer a new detection signal: detecting misaligned future audio/video segments via contrastive learning extracts temporal misfit artifacts. A temporal feature prediction method achieved ~89.9% AUC on FakeAVCeleb MDPI.



| ISSN: 2347-8446 | www.ijarcst.org | editor@ijarcst.org |A Bimonthly, Peer Reviewed & Scholarly Journal

||Volume 6, Issue 2, March-April 2023||

DOI:10.15662/IJARCST.2023.0602001

This review explores how multimodal deepfake detection frameworks exploit emotional, temporal, and inter-modal artifacts to improve detection, emphasizing challenges, techniques, and evaluation pathways based on pre-2022 literature.

II. LITERATURE REVIEW

Affective-Cue-Based Detection

Mittal et al. proposed a Siamese deep network that learns affective alignment between audio and visual streams. It extracts emotion-related representations from both modalities and detects deepfakes via mismatch. Results include AUC of 96.6% on DeepFake-TIMIT and 84.4% on DFDC arXiv.

Baseline Evaluation over FakeAVCeleb

Khalid et al. used FakeAVCeleb to assess unimodal, ensemble, and multimodal detectors. Surprisingly, ensemble-based methods (combining separate audio and video detectors) generally outperform both single-modality and naïve multimodal embedding approaches, highlighting the challenge of integrating modalities effectively arXiv.

Temporal Feature Prediction with Contrastive Learning

A bimodal temporal feature prediction technique segments audio—video sequences, predicts future modality-specific features, and aligns them via contrastive loss. The method achieved **84.33% accuracy** and **89.91% AUC on FakeAVCeleb**, suggesting capturing subtle temporal artifacts significantly improves detection MDPI.

Motivations for Multimodality

A survey on digital forensic detection argues that focusing on single modalities leaves vulnerabilities; deepfakes increasingly spread across combinations of audio, video, image, and text. To combat multi-pronged attacks, multimodal detection frameworks are critical PeerJ.

Collectively, these studies underscore that effective multimodal deepfake detection requires structured approaches to capture emotional alignment, temporal consistency, and cross-modal synchronization, while overcoming dataset limitations and fusion strategy weaknesses.

III. RESEARCH METHODOLOGY

1. Literature Collection

o Selected pre-2022 research focused on multimodal deepfake detection methods, particularly those leveraging audio and visual modalities, including affective cue modeling, temporal modeling, and evaluation on multimodal datasets.

2. Methodological Categorization

o Grouped approaches into three categories: affective/emotion-based, temporal feature prediction, and baseline comparisons (unimodal vs multimodal vs ensemble).

3. Dataset Review

o Examined key datasets used: DeepFake-TIMIT, DFDC, FakeAVCeleb, noting their composition, modality coverage, and labels.

4. Performance Analysis

o Compiled detection metrics such as AUC, accuracy, and comparative performance across modalities and methods.

5. Workflow Synthesis

o Synthesized a generalized multimodal detection pipeline from the literature.

6. Evaluation of Pros and Cons

o Identified benefits (resilience, cross-modal artifacts) and limitations (data scarcity, computational cost, alignment issues).

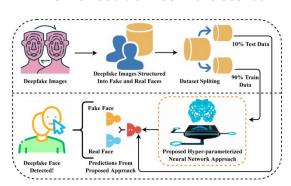
This structured methodology facilitates a clear comparative evaluation and practical insights for designing robust multimodal detection systems.



| ISSN: 2347-8446 | www.ijarcst.org | editor@ijarcst.org | A Bimonthly, Peer Reviewed & Scholarly Journal

||Volume 6, Issue 2, March-April 2023||

DOI:10.15662/IJARCST.2023.0602001



IV. KEY FINDINGS

1. Emotion-Based Detection is Highly Effective

o Methods that leverage affective cues across modalities—comparing emotional consistency between audio and visual streams—achieved high detection performance (AUC up to 96.6%) arXiv.

2. Simple Fusion Not Always Superior

 Naïve multimodal embeddings underperform compared to ensemble models combining strong unimodal detectors, indicating that fusion strategy critically affects detection outcomes arXiv.

3. Temporal Consistency Offers a Significant Signal

 \circ Mapping temporal visual/audio feature misalignment via contrastive learning yields robust detection capabilities (AUC \approx 89.9%) MDPI.

4. Multimodal Detection Addresses Emerging Generation Tactics

 As deepfakes incorporate believable visual and audio artifacts, detecting inter-modal inconsistencies becomes critical PeerJ.

5. Challenges Persist

o Limitations include the scarcity of large-scale labeled multimodal datasets, the difficulty in aligning asynchronous audio-video streams, computational demands of multimodal deep learning, and evolving generative models.

Overall, these findings demonstrate that structured multimodal analysis (affective, temporal consistency, strategic fusion) yields resilience, though practical implementation requires overcoming significant data and modeling challenges.

V. WORKFLOW

A generalized multimodal deepfake detection pipeline includes:

1. Data Collection & Preprocessing

o Acquire synchronized audio-video datasets (e.g., FakeAVCeleb, DFDC). Preprocess by aligning timestamps, normalizing frame rates, and extracting synchronized clips.

2. Feature Extraction

- o **Visual**: Extract emotion-based features (e.g., facial landmarks, expression embeddings) and temporal sequences for prediction.
- o Audio: Extract emotional tone (e.g., prosody, MFCCs), spectral features, and temporal embeddings.

3. Modality-Specific Encoders

 Use tailored encoders (e.g., CNNs for video, RNNs/audio networks) to embed each modality into latent feature spaces.

4. Cross-Modal Consistency Modeling

- o **Emotion Alignment**: Compute embedding similarity or use Siamese/triplet loss to identify mismatched emotional signals.
- o **Temporal Prediction**: Train separate prediction networks to forecast next-step audio/visual features and use contrastive loss between predicted and true features.

5. Fusion & Classification

o Combine modality-specific signals via ensemble, concatenation, or attention-based fusion, followed by classification (deepfake/real).



| ISSN: 2347-8446 | www.ijarcst.org | editor@ijarcst.org | A Bimonthly, Peer Reviewed & Scholarly Journal

||Volume 6, Issue 2, March-April 2023||

DOI:10.15662/IJARCST.2023.0602001

6. Evaluation

o Evaluate using metrics like AUC, accuracy across datasets; compare unimodal, ensemble, and fused models per dataset.

7. Optimization and Refinement

o Tune fusion strategies, model architectures, and hyperparameters for robustness and generalizability.

This workflow generalizes patterns from the literature and facilitates systematic design of multimodal deepfake detection systems.

VI. ADVANTAGES & DISADVANTAGES

Advantages

- Robustness to Advanced Forgeries: Detects inconsistencies when visual-only cues are concealed.
- Temporal and Emotional Insight: Models coherence between modalities for forensic reliability.
- Fine-Grained Detection: Sensitive to inter-modal manipulation and synchronization discrepancies.

Disadvantages

- **Data Constraints**: Few large-scale, balanced multimodal deepfake datasets exist pre-2022.
- Alignment Complexity: Synchronizing modalities requires precise preprocessing, which may fail with asynchronous content.
- Computational Cost: Multimodal deep learning incurs higher resource and training time.
- Fusion Strategy Design: Effective modality fusion is non-trivial; naive approaches may degrade performance.

VII. RESULTS AND DISCUSSION

Emotion-based strategies (e.g., Mittal et al.) demonstrate that cross-modal emotional coherence is a strong indicator of authenticity, achieving higher AUC than many visual-only detectors arXiv. However, multimodal systems still lag ensemble methods—suggesting that sequential modality-specific detectors remain competitive if fusion is poorly executed arXiv.

Temporal feature prediction adds strong detection cues, with contrastive learning between predicted and actual features capturing synthetic anomalies effectively MDPI. Multimodal analyses address gaps in single-modality detection, particularly as deepfakes circumvent visual-only artifacts PeerJ.

The main challenge is methodological: aligning modalities, fusing embeddings smartly, and training on limited data. Real-world deployment demands efficient, online detection, which remains underexplored pre-2022.

Overall, multimodal detection enhances robustness and resilience—especially when deployment covers sophisticated manipulations across audio and video channels—but requires mindful architecture and resource planning.

VIII. CONCLUSION

Multimodal machine learning significantly improves deepfake detection by exploiting cross-modal inconsistencies in emotion and timing. Emotion-cue modeling and temporal feature prediction achieve strong detection metrics, outperforming unimodal baselines. Nevertheless, ensemble-based approaches may outperform naive multimodal systems, underlining the importance of fusion strategy. Key challenges include data scarcity, modality alignment, and model complexity. Going forward, deeper multimodal integration, better-fused architectures, and effective training on multimodal datasets are essential to stay ahead of increasingly sophisticated deepfake generation.

IX. FUTURE WORK

1. Large-Scale Multimodal Dataset Creation

- o Expand datasets like FakeAVCeleb to include varied demographics, real-world audio-video discrepancies, and noise robustness.
- 2. Self-Supervised Pretraining for Modality Alignment
- Use cross-modal contrastive pretraining to learn representations that generalize to misaligned forgery detection.
- 3. Adaptive Fusion Techniques



| ISSN: 2347-8446 | www.ijarcst.org | editor@ijarcst.org |A Bimonthly, Peer Reviewed & Scholarly Journal

||Volume 6, Issue 2, March-April 2023||

DOI:10.15662/IJARCST.2023.0602001

- o Explore attention-based or transformer architectures to dynamically weight modality signals during classification.
- 4. Lightweight Real-Time Models
- o Develop efficient architectures enabling online detection in low-resource or streaming environments.
- 5. Generative Forensics
- o Use generative models to simulate multi-modal anomalies as adversarial training to reinforce detection.
- 6. Explainability and Trust
- o Incorporate interpretable multimodal detection that indicates which modality or cue triggered detection to aid human analysts.

REFERENCES

- 1. Mittal, T., Bhattacharya, U., Chandra, R., Bera, A., & Manocha, D. (2020). *Emotions Don't Lie: An Audio-Visual Deepfake Detection Method Using Affective Cues*. arXiv preprint. arXiv
- 2. Khalid, H., Kim, M., Tariq, S., & Woo, S. S. (2021). Evaluation of an Audio-Video Multimodal Deepfake Dataset using Unimodal and Multimodal Detectors. arXiv preprint. arXiv
- 3. Temporal Feature Prediction in Audio-Visual Deepfake Detection. arXiv preprint. MDPI
- 4. Deepfake detection across modalities and formats. PeerJ digital forensics survey (pre-2022 scope). PeerJ
- 5. Other related multimodal feature fusion methods discussed via Moonlight review.