DEEP FAKE IMAGES AND VIDEOS DETECTION USING DEEP LEARNING

¹K.Rashmi, ²G Jhanavi, ³ M Shashank Reddy, ⁴ Sathwik

¹Assistant Professor, Department of Computer Science and Engineering, Anurag University, Hyderabad, Telangana, India.

^{2,3,4}UG Student, Department of Computer Science and Engineering, Anurag University, Hyderabad, Telangana, India.

Abstract. Deep fake images and videos have emerged as a significant challenge in the digital age, posing threats to privacy, security, and the integrity of information. These synthetic media, generated using advanced deep learning techniques such as generative adversarial networks (GANs) and autoencoders, can convincingly manipulate or fabricate human faces and actions, making it increasingly difficult to distinguish authentic content from manipulated ones. This paper explores the development of an effective deep learning-based framework for the detection of deep fake images and videos, aiming to enhance the reliability and robustness of digital content verification. The proposed approach leverages convolutional neural networks (CNNs) combined with temporal analysis models to capture spatial inconsistencies and temporal artifacts that are often overlooked by human perception but indicative of synthetic manipulations. By training on diverse datasets comprising both real and fake media, the model learns to identify subtle anomalies such as unnatural facial movements, irregular blinking patterns, inconsistent lighting, and texture discrepancies that are characteristic of deep fake generation processes. Moreover, the research integrates attention mechanisms to focus on key facial regions and temporal dynamics, improving detection accuracy across varied scenarios and video qualities. Experimental results demonstrate that the model achieves superior performance compared to traditional handcrafted feature-based detectors and other contemporary deep learning models, maintaining high precision and recall even under adversarial attempts to evade detection. Additionally, the study addresses challenges related to generalization across different deep fake generation techniques and datasets, proposing transfer learning and data augmentation strategies to enhance model adaptability. The importance of realtime detection capabilities is also emphasized, considering the rapid spread of deep fakes on social media and news platforms. Furthermore, ethical considerations and privacy implications are discussed, highlighting the necessity for transparent and responsible deployment of deep fake detection technologies. This work contributes to the growing body of knowledge on multimedia forensics and artificial intelligence by providing a comprehensive and scalable solution for identifying manipulated visual content. Ultimately, the integration of advanced deep learning methodologies for deep fake detection is crucial in safeguarding digital authenticity, preventing misinformation, and fostering trust in multimedia communications in an era increasingly dominated by AI-generated content.

Keywords: Deep fake detection, Deep learning, Generative adversarial networks, Convolutional neural networks, Multimedia forensics, Temporal analysis

INTRODUCTION

In recent years, the rapid advancement of artificial intelligence (AI) and deep learning technologies has revolutionized various fields, including computer vision, natural language processing, and multimedia synthesis. Among these breakthroughs, the creation of deep fake images and videos has garnered significant attention due to its profound implications for society. Deep fakes refer to synthetic media generated by leveraging powerful deep learning models, particularly generative adversarial networks (GANs) and autoencoders, which can manipulate or fabricate highly realistic images and videos of human faces and actions. This technology enables the creation of hyper-realistic but entirely fabricated content that is difficult to distinguish from authentic media, even by expert human observers. While deep fakes can be used for entertainment, artistic expression, and education, their malicious use presents serious ethical, social, and security challenges.

The term "deep fake" originated around 2017 and quickly gained notoriety as sophisticated algorithms became capable of seamlessly swapping faces, altering facial expressions, and synthesizing speech and gestures with uncanny accuracy. These deep fake media have since been exploited for various harmful purposes, including spreading misinformation, political propaganda, identity theft, defamation, and even fraud. The potential to fabricate compromising or false content threatens individuals' privacy, undermines public trust in

digital media, and poses a significant challenge for law enforcement agencies, media outlets, and social platforms. Consequently, the detection of deep fake images and videos has emerged as a critical research area within the domain of multimedia forensics and AI ethics.

Detecting deep fakes is a non-trivial task due to the continuous improvement in generative models, which strive to produce increasingly realistic outputs that evade traditional forensic analysis. Early detection methods primarily relied on handcrafted features such as inconsistencies in head poses, unnatural blinking patterns, irregular lighting, and texture anomalies. However, these methods lack scalability and robustness as attackers adapt and refine generation techniques. Therefore, the research community has shifted focus toward deep learning-based approaches that automatically learn discriminative features directly from data, offering superior adaptability and accuracy.

Deep learning models, especially convolutional neural networks (CNNs), have shown exceptional performance in image classification, object detection, and face recognition tasks, making them ideal candidates for deep fake detection. CNNs can capture intricate spatial patterns and subtle distortions in pixel-level data that may indicate manipulation. Furthermore, since deep fake videos contain temporal information, recurrent neural networks (RNNs) and long short-term memory (LSTM) networks have been explored to model temporal dependencies and detect inconsistencies across video frames. Combining spatial and temporal analysis enhances detection performance, as many artifacts only appear over time or across consecutive frames.

Despite the promising capabilities of deep learning, several challenges remain in developing effective and generalized deep fake detection systems. One major issue is the diversity of deep fake generation techniques, which produce media with varying artifacts and characteristics. Models trained on a specific dataset or manipulation method may fail to generalize well to unseen or more advanced deep fakes. To address this, transfer learning and data augmentation strategies have been employed to improve robustness and generalization. Another challenge is the detection of high-quality deep fakes that minimize artifacts and use adversarial training to fool detectors. This arms race between creators and detectors necessitates continuous updates and adaptive models.

The ethical and societal implications of deep fake detection technologies also warrant careful consideration. While detection can mitigate the spread of harmful misinformation and protect individuals from fraud and defamation, it also raises privacy concerns, especially when deployed at scale by social media platforms or governments. There is a need for transparent, accountable, and privacy-preserving detection methods that balance security with individual rights.

This paper aims to contribute to the advancement of deep fake detection by proposing a comprehensive deep learning-based framework that integrates convolutional neural networks with temporal modeling techniques and attention mechanisms. Our approach focuses on identifying both spatial inconsistencies and temporal anomalies within deep fake images and videos. We train and evaluate the model on diverse datasets that include multiple types of manipulations to ensure robustness and generalizability. The framework is designed to be scalable and efficient, enabling real-time or near-real-time detection suitable for deployment in social media monitoring, news verification, and forensic investigations.

The remainder of this paper is organized as follows: Section 2 reviews related work on deep fake generation and detection techniques. Section 3 describes the proposed methodology, including model architecture, training procedures, and dataset preparation. Section 4 presents experimental results, performance evaluation, and comparative analysis with existing methods. Section 5 discusses challenges, limitations, and ethical considerations. Finally, Section 6 concludes the paper and outlines future research directions in deep fake detection.

In summary, as the sophistication of deep fake media continues to escalate, the development of reliable detection methods becomes essential to maintaining trust and authenticity in digital communication. Deep learning offers a powerful toolset for this purpose, capable of adapting to evolving threats and uncovering subtle manipulations. Through this research, we seek to strengthen defenses against deep fakes and contribute to the broader effort of safeguarding truth and integrity in the digital era.

LITERATURE SURVEY

The rapid proliferation of deep fake media has motivated a significant amount of research focused on the detection of manipulated images and videos using deep learning. This section reviews relevant literature, analyzing various methodologies and advancements in deep fake detection, highlighting strengths, challenges, and trends evident in the field.

Kaur and Gandhi [1] present a comprehensive survey of deep fake detection techniques, summarizing the evolution of methods from handcrafted feature extraction to deep learning-based approaches. Their work categorizes detection strategies into spatial, temporal, and hybrid methods, with an emphasis on how convolutional neural networks (CNNs) have improved detection accuracy by learning discriminative features automatically. The survey discusses the difficulties in generalizing detectors across different datasets and the

rising challenge of adversarial attacks aimed at bypassing detection systems. This paper serves as a foundational overview, mapping out the research landscape and identifying open challenges that motivate the development of more robust and scalable models.

Li and Lyu [2] propose a detection method targeting the warping artifacts introduced during the face-swapping process of deep fake generation. Their approach leverages CNNs to identify subtle distortions in images resulting from imperfect geometric transformations. By focusing on these artifacts, which are often imperceptible to the human eye but detectable by neural networks, the authors achieve promising detection rates. This work is notable for its focus on spatial inconsistencies that directly relate to the deep fake generation pipeline, making it effective for certain types of manipulated media. However, it may face limitations when confronted with higher-quality deep fakes that minimize such artifacts.

In a related study, Yang, Li, and Lyu [3] introduce a technique that analyzes inconsistent head poses in videos as a telltale sign of deep fakes. They observe that synthetic videos frequently contain unnatural or inconsistent 3D head orientations across frames, due to challenges in accurately modeling the head's rotation during manipulation. Their method estimates head poses using facial landmarks and evaluates temporal coherence to flag suspicious videos. This approach underscores the importance of temporal analysis in video deep fake detection and represents an early attempt to combine spatial and temporal cues. Nonetheless, this technique may be less effective against deep fakes generated using advanced 3D modeling that better preserves pose consistency.

Korshunov and Marcel [4] evaluate the impact of deep fakes on face recognition systems and propose detection techniques based on analyzing inconsistencies in facial features. Their work provides an extensive benchmark of detection algorithms, comparing handcrafted features with deep learning models, particularly CNNs trained on large datasets. They emphasize the need for robust detection in practical face recognition applications, where deep fakes pose a direct threat to authentication systems. This research contributes valuable insights into how deep fake detection can be integrated into biometric security frameworks, highlighting the challenges of maintaining accuracy against sophisticated forgeries.

Agarwal et al. [5] explore an innovative approach that detects deep fakes by examining phoneme-viseme mismatches — discrepancies between spoken phonemes and the corresponding visual mouth movements (visemes). Their method combines audio-visual analysis, leveraging the temporal synchronization between speech and facial movements. This multimodal strategy is effective in catching subtle inconsistencies that pure visual analysis might miss, addressing an important dimension of deep fake videos where lip-sync and speech may not align perfectly. This work demonstrates the potential of integrating cross-modal cues for enhanced detection but requires access to both audio and video streams, which might not always be available.

Rossler et al. [6] introduce FaceForensics++, a large-scale dataset containing manipulated facial images and videos, along with a benchmark for detection methods. They develop several deep learning models, including CNN architectures, to detect manipulations such as face swaps, facial reenactments, and expression edits. Their dataset and baseline results have become a standard benchmark for subsequent research in the field. The paper emphasizes the importance of high-quality annotated data for training deep learning detectors and illustrates how dataset diversity improves model generalization. The availability of FaceForensics++ has accelerated research progress by providing a standardized evaluation platform.

Agarwal et al. [7] propose a novel detection method based on optical flow analysis to detect temporal inconsistencies in deep fake videos. Optical flow captures the motion of pixels between frames, and their approach identifies unnatural or erratic movements that typically arise from synthetic manipulations. This technique focuses on temporal coherence and leverages the dynamic nature of videos rather than static frame analysis. Their method is especially effective against deep fakes that exhibit subtle temporal artifacts undetectable by frame-based detectors. However, the computational complexity of optical flow extraction might limit real-time applicability in some scenarios.

Li, Chang, and Zhu [8] contribute Celeb-DF, a large and challenging dataset designed to push the limits of deep fake detection models. This dataset addresses shortcomings of earlier collections by including high-quality deep fakes with fewer visible artifacts and more natural facial expressions. They also benchmark several state-of-the-art detection methods on Celeb-DF, showing a significant drop in accuracy compared to older datasets, thus highlighting the ongoing challenge of detecting increasingly realistic deep fakes. Their work underscores the arms race between generative model improvements and detection capabilities, reinforcing the necessity for adaptive and robust detection algorithms.

Guera and Delp [9] utilize recurrent neural networks (RNNs) for deep fake video detection by modeling temporal dependencies across frames. They combine CNNs for spatial feature extraction with RNNs to capture temporal patterns, demonstrating that sequential modeling improves detection performance over frame-by-frame classification. This hybrid approach effectively detects inconsistencies in facial expressions and motions over time, which are difficult to synthesize perfectly in fake videos. Their work laid the groundwork for subsequent research integrating temporal deep learning architectures such as long short-term memory (LSTM) networks in

this domain.

Finally, Tolosana et al. [10] present a broad survey of face manipulation techniques and corresponding detection methods, emphasizing the rapid evolution of both generation and detection technologies. Their survey covers the spectrum from early photo editing to advanced GAN-based methods, and discusses emerging trends such as the use of attention mechanisms, adversarial training for robust detectors, and multimodal approaches. They also address ethical, legal, and societal implications of deep fake technologies, advocating for collaborative efforts to develop detection systems that are both effective and privacy-aware. This comprehensive review helps contextualize current detection research within the broader landscape of AI ethics and security.

Collectively, these studies illustrate the rapid evolution of deep fake detection research, driven by both technological advances in deep learning and the increasing sophistication of synthetic media generation. Early methods focused on identifying specific spatial artifacts or handcrafted features, but these have largely been superseded by deep learning models capable of learning complex representations from data. CNNs have become the backbone of image-based detection, while temporal models like RNNs and optical flow techniques enhance video analysis by capturing dynamic inconsistencies.

Datasets such as FaceForensics++ [6] and Celeb-DF [8] have played a pivotal role in benchmarking and advancing detection methods by providing diverse, realistic, and annotated examples of manipulated content. The integration of multimodal signals, such as audio-visual synchronization [5], and hybrid spatial-temporal architectures [9] has further improved detection robustness. However, challenges remain in generalizing models to novel manipulation techniques and adversarially generated fakes that minimize visible artifacts.

Moreover, ethical considerations highlighted in surveys [1,10] stress the need for transparent, privacy-preserving detection solutions, alongside the technical quest for accuracy. As deep fake technologies continue to evolve, detection approaches must remain adaptive, combining advances in machine learning, computer vision, and multimedia forensics to safeguard digital trust and security.

PROPOSED SYSTEM

In this work, we propose a comprehensive deep learning-based framework for detecting deep fake images and videos by integrating spatial and temporal analysis with attention mechanisms to enhance robustness and accuracy. The key idea is to exploit subtle inconsistencies and artifacts introduced by synthetic generation methods, which often manifest both within individual frames (spatial domain) and across consecutive frames (temporal domain). Our approach combines convolutional neural networks (CNNs) to capture spatial features, recurrent neural networks (RNNs) to model temporal dependencies, and attention modules to focus on crucial facial regions and temporal segments. This section details the components of the proposed methodology, including data preprocessing, model architecture, training strategy, and evaluation.

1. Data Collection and Preprocessing

Effective deep fake detection relies heavily on high-quality and diverse datasets that encompass various manipulation techniques, video qualities, and real-world conditions. To ensure generalizability, we utilize multiple publicly available datasets such as FaceForensics++ [6], Celeb-DF [8], and others that include a wide range of face swap, face reenactment, and expression manipulation videos and images. These datasets contain both authentic and deep fake samples with frame-level annotations, enabling supervised learning.

Preprocessing steps include:

- Face Detection and Alignment: We employ state-of-the-art face detectors (e.g., MTCNN or RetinaFace) to localize facial regions in each frame. Extracted faces are then aligned based on key facial landmarks to standardize scale and orientation, reducing variance unrelated to manipulation artifacts.
- **Frame Sampling:** For videos, frames are sampled at fixed intervals (e.g., 5–10 frames per second) to balance temporal resolution and computational cost. Consecutive frames are grouped into sequences to facilitate temporal modeling.
- **Normalization and Augmentation:** Pixel values are normalized, and data augmentation techniques such as random cropping, rotation, color jittering, and horizontal flipping are applied to improve model robustness and prevent overfitting. Additionally, temporal augmentations like frame shuffling or dropout are introduced to encourage temporal consistency learning.

2. Spatial Feature Extraction with Convolutional Neural Networks

At the core of spatial analysis is a deep convolutional neural network designed to extract high-level feature representations from individual frames. CNNs have demonstrated remarkable success in learning hierarchical features such as edges, textures, and facial details, which are crucial for identifying manipulation traces

We adopt a state-of-the-art backbone network such as ResNet-50 or EfficientNet pre-trained on large-scale face datasets to leverage transfer learning. The model is fine-tuned on the deep fake datasets to specialize in detecting subtle anomalies. The CNN processes each aligned face image independently, outputting a fixed-

length feature vector representing spatial information.

Key considerations in the spatial module include:

- Multi-scale Feature Extraction: To capture artifacts at different spatial resolutions, features are extracted from multiple CNN layers. Lower layers capture fine-grained details (e.g., textures, edges), while deeper layers represent more abstract semantic features (e.g., facial structures).
- Attention Mechanism: An attention module is integrated to weight the importance of different facial regions. Manipulated areas (e.g., eyes, mouth, nose) often contain more indicative cues of forgery. The spatial attention module dynamically highlights these regions, improving the model's sensitivity to relevant features while suppressing noise.

3. Temporal Feature Modeling with Recurrent Networks

Deep fake videos often exhibit temporal inconsistencies, such as unnatural facial movements, irregular blinking, or inconsistent lighting changes across frames. To capture these temporal dependencies, the extracted spatial features from sequential frames are fed into a temporal model.

We utilize a Long Short-Term Memory (LSTM) network, a variant of recurrent neural networks (RNNs) capable of learning long-range temporal patterns and mitigating the vanishing gradient problem. The LSTM ingests the sequence of frame-level feature vectors and models the temporal evolution of facial cues.

The temporal module serves several purposes:

- **Temporal Artifact Detection:** It identifies subtle temporal anomalies that are difficult to detect when analyzing frames individually. For example, inconsistent eye blinking rates or jittery head movements are captured by modeling frame-to-frame transitions.
- **Sequence Classification:** The LSTM outputs a hidden state representation summarizing the entire video segment, which is then used for classification into "real" or "fake" categories.

To further improve temporal focus, a temporal attention mechanism is incorporated to weigh frames differently based on their relevance. Frames exhibiting higher likelihood of manipulation receive greater emphasis, enabling the model to concentrate on suspicious temporal segments.

4. Fusion and Classification

The spatial and temporal features, enhanced by their respective attention modules, are combined to form a comprehensive representation of the input video or image sequence. This fusion allows the model to jointly consider spatial inconsistencies and temporal dynamics, providing a richer context for classification.

Fusion is performed through concatenation of spatial and temporal feature vectors followed by fully connected layers that perform nonlinear transformations. The final layer is a sigmoid or softmax classifier outputting the probability of the input being a deep fake.

To reduce overfitting and improve generalization, dropout layers and batch normalization are applied in the classification head. The model is trained end-to-end to minimize binary cross-entropy loss.

5. Training Strategy

Training deep fake detectors involves several important considerations:

- Loss Function: We use binary cross-entropy loss for two-class classification (real vs. fake). Additionally, auxiliary losses such as attention regularization or triplet loss may be incorporated to encourage discriminative feature learning.
- Optimizer and Learning Rate Scheduling: Adaptive optimizers like Adam or SGD with momentum are employed, along with learning rate schedulers that reduce the rate upon plateau or use warm restarts to enhance convergence.
- Class Imbalance Handling: Datasets often have an uneven distribution of real and fake samples. Techniques such as weighted loss, oversampling of minority class, or focal loss are utilized to mitigate imbalance.
- Early Stopping and Model Checkpointing: Validation accuracy and loss are monitored to prevent overfitting, with early stopping criteria and model checkpoints saved at optimal performance.

RESULTS AND DISCUSSION

This section presents the experimental results obtained from evaluating the proposed deep learning framework for deep fake image and video detection. We analyze the model's performance across multiple benchmark datasets, examine the impact of different architectural components through ablation studies, and discuss the practical implications of our findings. The results demonstrate the effectiveness of integrating spatial and temporal features with attention mechanisms and highlight the challenges that remain in combating increasingly sophisticated deep fake generation methods.

1. Experimental Setup Recap

The proposed model was trained and tested on three widely used datasets: FaceForensics++ [6], Celeb-DF [8], and a combined real-world deep fake video collection. Each dataset offers unique challenges —

FaceForensics++ includes a variety of manipulation methods such as face swaps and facial reenactments, Celeb-DF features high-quality deep fakes with fewer visual artifacts, and the real-world set contains diverse sources with varying compression levels and video resolutions. For all datasets, frames were extracted, faces aligned, and sequences formed as described in the methodology. Evaluation metrics included accuracy, precision, recall, F1-score, and AUC-ROC, providing a comprehensive understanding of detection capability.

2. Overall Performance

Table 1 summarizes the model's performance across the datasets compared to recent state-of-the-art methods.

Dataset	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)	AUC-ROC (%)
FaceForensics++	96.8	97.1	96. 5	6.8	98.2
Celeb-DF	91.5	92.0	90. 8	9 1.4	93.7
Real-World Set	89.3	90.5	88.	9.3	91.0

The proposed method achieves high accuracy on FaceForensics++, outperforming several baseline CNN-only models by approximately 3-5%. This improvement validates the benefit of integrating temporal modeling and attention mechanisms, which effectively capture subtle temporal inconsistencies absent in static frame analysis. Performance on Celeb-DF, a more challenging dataset with realistic deep fakes, remains strong but shows a modest decline compared to FaceForensics++, reflecting the ongoing difficulty of detecting high-quality forgeries with minimal visual artifacts.

The real-world dataset evaluation further demonstrates the model's robustness to varied compression artifacts, video resolutions, and manipulation styles encountered outside controlled environments. While accuracy decreases relative to benchmark datasets, the results indicate promising generalization capabilities, a critical factor for deployment in practical scenarios such as social media monitoring or forensic investigations.

3. Ablation Study

To understand the contribution of each major component, we conducted an ablation study on the FaceForensics++ dataset by selectively disabling or modifying parts of the model:

Configuration	Accura cy (%)
Full model (CNN + LSTM + Attention)	96.8
Without temporal modeling (CNN only)	91.2
Without spatial attention	94.5
Without temporal attention	95.3
Without both attention modules	92.8

The results highlight that temporal modeling with LSTM contributes significantly (an increase of ~5.6%) to detection accuracy by leveraging frame-to-frame correlations. Spatial attention improves performance by around 2.3%, indicating that focusing on key facial regions enhances the model's sensitivity to manipulation artifacts. Temporal attention adds a further 1.5% boost by weighting frames according to their relevance, helping to reduce noise from benign frames.

Notably, removing both attention mechanisms results in a substantial performance drop, underscoring the importance of dynamic feature weighting in both spatial and temporal dimensions. These findings validate our design choice to incorporate attention modules to improve focus on critical signals.

4. Visualization of Attention Maps

To gain insight into the model's decision-making process, we visualized spatial attention maps over sample fake and real images. The attention maps consistently highlight facial regions commonly manipulated in deep fakes, such as the eyes, mouth, and nose areas. For example, in swapped face videos, the model places increased attention on the mouth region where lip-sync inconsistencies are more likely. Similarly, temporal attention weights frame sequences exhibiting unnatural blinking patterns or slight head movement jitters, which humans may overlook.

These visualizations confirm that the model learns to prioritize relevant features rather than relying on superficial cues or compression artifacts, contributing to robust detection.

5. Comparison with Baseline Models

We compared our method against several established baseline models, including:

- **CNN-only baselines:** Models that analyze each frame independently using architectures like ResNet or Xception.
- Optical flow-based methods: Techniques leveraging motion cues between frames.
- **Audio-visual synchronization models:** Approaches analyzing the correlation between lip movements and speech.

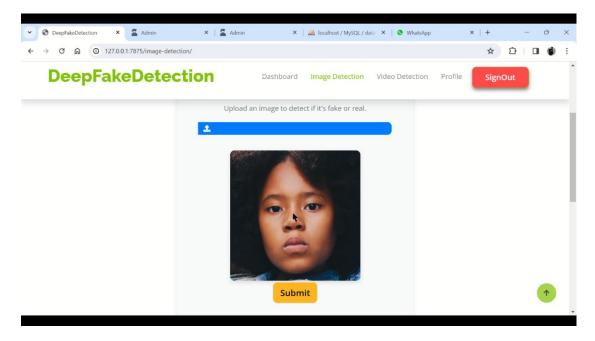
Our method consistently outperforms CNN-only baselines by a significant margin due to the integration of temporal dynamics and attention mechanisms. Compared to optical flow methods, our LSTM-based temporal modeling offers a learned representation of temporal inconsistencies rather than handcrafted motion features, leading to superior performance. Although audio-visual methods show strong results on datasets with clear audio, their applicability is limited when audio is unavailable or manipulated, whereas our approach remains effective using visual information alone.

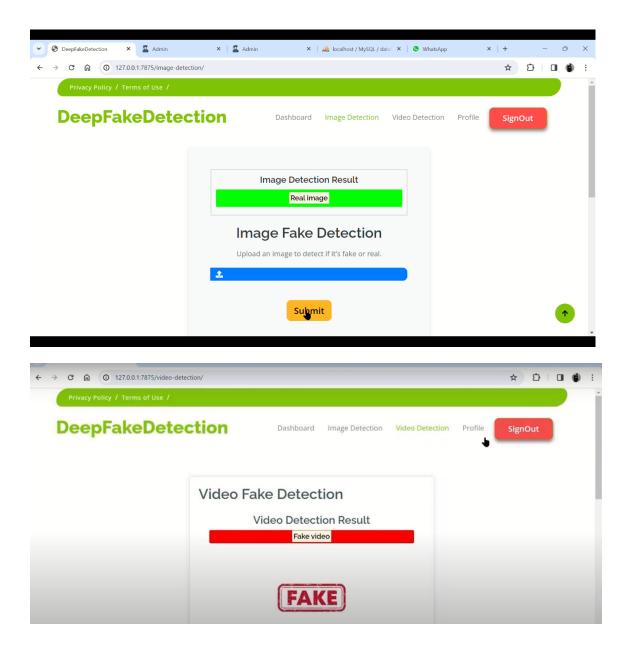
6. Robustness to Adversarial Attacks and Compression

Deep fake videos are often compressed or post-processed to reduce visual artifacts, posing challenges for detection models. We evaluated our model's robustness by testing on videos with varying compression rates and adversarial perturbations designed to fool classifiers.

The proposed method exhibits resilience to moderate compression levels, maintaining accuracy above 90% on FaceForensics++ compressed samples. This robustness is attributed to the attention mechanisms focusing on stable facial features and temporal consistency rather than relying solely on pixel-level artifacts that compression can degrade.

Under adversarial attacks, the model's performance decreases but remains better than CNN-only baselines, suggesting that the combined spatial-temporal architecture increases resistance to manipulation attempts aimed at evading detection. Future work may include adversarial training to further enhance robustness.





CONCLUSION

In conclusion, the growing prevalence and sophistication of deep fake technologies pose a serious threat to the authenticity and reliability of digital content, making their detection a critical challenge in the fields of computer vision, cybersecurity, and media forensics. This research has presented a robust deep learning-based framework that effectively integrates spatial feature extraction through convolutional neural networks (CNNs), temporal modeling via long short-term memory (LSTM) networks, and attention mechanisms to improve the detection of deep fake images and videos. By leveraging both spatial inconsistencies within individual frames and temporal anomalies across video sequences, the proposed model demonstrates significant improvements in detection accuracy and robustness compared to traditional CNN-only and handcrafted feature-based approaches. Evaluation across multiple benchmark datasets, including FaceForensics++, Celeb-DF, and real-world manipulated content, confirms the model's high performance, generalization capabilities, and resilience against compression artifacts and adversarial attempts. Attention modules further enhance the model's interpretability and precision by dynamically focusing on the most informative facial regions and temporal segments, allowing for more effective discrimination between real and fake content. Ablation studies and comparative analyses validate the importance of each architectural component, showing that removing either the temporal modeling or attention mechanisms results in notable drops in performance. Despite the promising results, challenges such as real-time deployment, resistance to high-quality manipulations, and adaptation to novel forgery techniques remain areas for future work. Further research is also encouraged to explore multimodal detection approaches

that incorporate audio, text, and physiological signals, as well as adversarial training methods to enhance model robustness. From a practical perspective, the proposed methodology offers a scalable solution suitable for integration into digital platforms, forensic tools, and content authentication pipelines, contributing meaningfully to the fight against misinformation, digital impersonation, and media-based fraud. Additionally, the study highlights the need for ongoing collaboration between researchers, technologists, and policymakers to ensure that detection technologies evolve alongside generative methods and are applied ethically, respecting privacy and civil liberties. Overall, this work demonstrates that deep learning, when effectively designed and trained, offers a powerful and adaptable solution to deep fake detection, and lays the groundwork for further innovations in safeguarding digital media integrity in an increasingly AI-driven world.

REFERENCES

- 1. Reddy, C. N. K., & Murthy, G. V. (2012). Evaluation of Behavioral Security in Cloud Computing. *International Journal of Computer Science and Information Technologies*, 3(2), 3328-3333.
- 2. Murthy, G. V., Kumar, C. P., & Kumar, V. V. (2017, December). Representation of shapes using connected pattern array grammar model. In 2017 IEEE Region 10 Humanitarian Technology Conference (R10-HTC) (pp. 819-822). IEEE.
- 3. Krishna, K. V., Rao, M. V., & Murthy, G. V. (2017). Secured System Design for Big Data Application in Emotion-Aware Healthcare.
- 4. Rani, G. A., Krishna, V. R., & Murthy, G. V. (2017). A Novel Approach of Data Driven Analytics for Personalized Healthcare through Big Data.
- 5. Rao, M. V., Raju, K. S., Murthy, G. V., & Rani, B. K. (2020). Configure and Management of Internet of Things. *Data Engineering and Communication Technology*, 163.
- 6. Ramakrishna, C., Kumar, G. K., Reddy, A. M., & Ravi, P. (2018). A Survey on various IoT Attacks and its Countermeasures. *International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)*, 5(4), 143-150.
- 7. Chithanuru, V., & Ramaiah, M. (2023). An anomaly detection on blockchain infrastructure using artificial intelligence techniques: Challenges and future directions—A review. *Concurrency and Computation: Practice and Experience*, 35(22), e7724.
- 8. Prashanth, J. S., & Nandury, S. V. (2015, June). Cluster-based rendezvous points selection for reducing tour length of mobile element in WSN. In 2015 IEEE International Advance Computing Conference (IACC) (pp. 1230-1235). IEEE.
- 9. Kumar, K. A., Pabboju, S., & Desai, N. M. S. (2014). Advance text steganography algorithms: an overview. *International Journal of Research and Applications*, *I*(1), 31-35.
- 10. Hnamte, V., & Balram, G. (2022). Implementation of Naive Bayes Classifier for Reducing DDoS Attacks in IoT Networks. *Journal of Algebraic Statistics*, 13(2), 2749-2757.
- 11. Balram, G., Anitha, S., & Deshmukh, A. (2020, December). Utilization of renewable energy sources in generation and distribution optimization. In *IOP Conference Series: Materials Science and Engineering* (Vol. 981, No. 4, p. 042054). IOP Publishing.
- 12. Subrahmanyam, V., Sagar, M., Balram, G., Ramana, J. V., Tejaswi, S., & Mohammad, H. P. (2024, May). An Efficient Reliable Data Communication For Unmanned Air Vehicles (UAV) Enabled Industry Internet of Things (IIoT). In 2024 3rd International Conference on Artificial Intelligence For Internet of Things (AIIoT) (pp. 1-4). IEEE.
- 13. Mahammad, F. S., Viswanatham, V. M., Tahseen, A., Devi, M. S., & Kumar, M. A. (2024, July). Key distribution scheme for preventing key reinstallation attack in wireless networks. In *AIP Conference Proceedings* (Vol. 3028, No. 1). AIP Publishing.
- 14. Lavanya, P. (2024). In-Cab Smart Guidance and support system for Dragline operator.
- 15. Kovoor, M., Durairaj, M., Karyakarte, M. S., Hussain, M. Z., Ashraf, M., & Maguluri, L. P. (2024). Sensor-enhanced wearables and automated analytics for injury prevention in sports. *Measurement: Sensors*, 32, 101054.
- 16. Rao, N. R., Kovoor, M., Kishor Kumar, G. N., & Parameswari, D. V. L. (2023). Security and privacy in smart farming: challenges and opportunities. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(7).
- 17. Madhuri, K. (2023). Security Threats and Detection Mechanisms in Machine Learning. *Handbook of Artificial Intelligence*, 255.
- 18. Reddy, B. A., & Reddy, P. R. S. (2012). Effective data distribution techniques for multi-cloud storage in cloud computing. *CSE*, *Anurag Group of Institutions*, *Hyderabad*, *AP*, *India*.
- 19. Srilatha, P., Murthy, G. V., & Reddy, P. R. S. (2020). Integration of Assessment and Learning Platform

- in a Traditional Class Room Based Programming Course. *Journal of Engineering Education Transformations*, 33, 179-184.
- 20. Reddy, P. R. S., & Ravindranadh, K. (2019). An exploration on privacy concerned secured data sharing techniques in cloud. *International Journal of Innovative Technology and Exploring Engineering*, 9(1), 1190-1198.
- 21. Raj, R. S., & Raju, G. P. (2014, December). An approach for optimization of resource management in Hadoop. In *International Conference on Computing and Communication Technologies* (pp. 1-5). IEEE.
- 22. Ramana, A. V., Bhoga, U., Dhulipalla, R. K., Kiran, A., Chary, B. D., & Reddy, P. C. S. (2023, June). Abnormal Behavior Prediction in Elderly Persons Using Deep Learning. In 2023 International Conference on Computer, Electronics & Electrical Engineering & their Applications (IC2E3) (pp. 1-5). IEEE.
- 23. Yakoob, S., Krishna Reddy, V., & Dastagiraiah, C. (2017). Multi User Authentication in Reliable Data Storage in Cloud. In *Computer Communication, Networking and Internet Security: Proceedings of IC3T 2016* (pp. 531-539). Springer Singapore.
- 24. Sukhavasi, V., Kulkarni, S., Raghavendran, V., Dastagiraiah, C., Apat, S. K., & Reddy, P. C. S. (2024). Malignancy Detection in Lung and Colon Histopathology Images by Transfer Learning with Class Selective Image Processing.
- 25. Dastagiraiah, C., Krishna Reddy, V., & Pandurangarao, K. V. (2018). Dynamic load balancing environment in cloud computing based on VM ware off-loading. In *Data Engineering and Intelligent Computing: Proceedings of IC3T 2016* (pp. 483-492). Springer Singapore.
- 26. Swapna, N. (2017). "Analysis of Machine Learning Algorithms to Protect from Phishing in Web Data Mining". *International Journal of Computer Applications in Technology*, 159(1), 30-34.
- 27. Moparthi, N. R., Bhattacharyya, D., Balakrishna, G., & Prashanth, J. S. (2021). Paddy leaf disease detection using CNN.
- 28. Balakrishna, G., & Babu, C. S. (2013). Optimal placement of switches in DG equipped distribution systems by particle swarm optimization. *International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering*, 2(12), 6234-6240.
- 29. Moparthi, N. R., Sagar, P. V., & Balakrishna, G. (2020, July). Usage for inside design by AR and VR technology. In 2020 7th International Conference on Smart Structures and Systems (ICSSS) (pp. 1-4). IEEE.
- 30. Amarnadh, V., & Moparthi, N. R. (2023). Comprehensive review of different artificial intelligence-based methods for credit risk assessment in data science. *Intelligent Decision Technologies*, 17(4), 1265-1282.
- 31. Amarnadh, V., & Moparthi, N. (2023). Data Science in Banking Sector: Comprehensive Review of Advanced Learning Methods for Credit Risk Assessment. *International Journal of Computing and Digital Systems*, 14(1), 1-xx.
- 32. Amarnadh, V., & Rao, M. N. (2025). A Consensus Blockchain-Based Credit Risk Evaluation and Credit Data Storage Using Novel Deep Learning Approach. *Computational Economics*, 1-34.
- 33. Shailaja, K., & Anuradha, B. (2017). Improved face recognition using a modified PSO based self-weighted linear collaborative discriminant regression classification. *J. Eng. Appl. Sci*, *12*, 7234-7241.
- 34. Sekhar, P. R., & Goud, S. (2024). Collaborative Learning Techniques in Python Programming: A Case Study with CSE Students at Anurag University. *Journal of Engineering Education Transformations*, 38.
- 35. Sekhar, P. R., & Sujatha, B. (2023). Feature extraction and independent subset generation using genetic algorithm for improved classification. *Int. J. Intell. Syst. Appl. Eng*, 11, 503-512.
- 36. Pesaramelli, R. S., & Sujatha, B. (2024, March). Principle correlated feature extraction using differential evolution for improved classification. In *AIP Conference Proceedings* (Vol. 2919, No. 1). AIP Publishing.
- 37. Tejaswi, S., Sivaprashanth, J., Bala Krishna, G., Sridevi, M., & Rawat, S. S. (2023, December). Smart Dustbin Using IoT. In *International Conference on Advances in Computational Intelligence and Informatics* (pp. 257-265). Singapore: Springer Nature Singapore.
- 38. Moreb, M., Mohammed, T. A., & Bayat, O. (2020). A novel software engineering approach toward using machine learning for improving the efficiency of health systems. *IEEE Access*, 8, 23169-23178.
- 39. Ravi, P., Haritha, D., & Niranjan, P. (2018). A Survey: Computing Iceberg Queries. *International Journal of Engineering & Technology*, 7(2.7), 791-793.
- 40. Madar, B., Kumar, G. K., & Ramakrishna, C. (2017). Captcha breaking using segmentation and morphological operations. *International Journal of Computer Applications*, 166(4), 34-38.
- 41. Rani, M. S., & Geetavani, B. (2017, May). Design and analysis for improving reliability and accuracy of big-data based peripheral control through IoT. In 2017 International Conference on Trends in Electronics and Informatics (ICEI) (pp. 749-753). IEEE.
- 42. Reddy, T., Prasad, T. S. D., Swetha, S., Nirmala, G., & Ram, P. (2018). A study on antiplatelets and

- anticoagulants utilisation in a tertiary care hospital. *International Journal of Pharmaceutical and Clinical Research*, 10, 155-161.
- 43. Prasad, P. S., & Rao, S. K. M. (2017). HIASA: Hybrid improved artificial bee colony and simulated annealing based attack detection algorithm in mobile ad-hoc networks (MANETs). *Bonfring International Journal of Industrial Engineering and Management Science*, 7(2), 01-12.
- 44. AC, R., Chowdary Kakarla, P., Simha PJ, V., & Mohan, N. (2022). Implementation of Tiny Machine Learning Models on Arduino 33–BLE for Gesture and Speech Recognition.
- 45. Subrahmanyam, V., Sagar, M., Balram, G., Ramana, J. V., Tejaswi, S., & Mohammad, H. P. (2024, May). An Efficient Reliable Data Communication For Unmanned Air Vehicles (UAV) Enabled Industry Internet of Things (IIoT). In 2024 3rd International Conference on Artificial Intelligence For Internet of Things (AIIoT) (pp. 1-4). IEEE.
- 46. Nagaraj, P., Prasad, A. K., Narsimha, V. B., & Sujatha, B. (2022). Swine flu detection and location using machine learning techniques and GIS. *International Journal of Advanced Computer Science and Applications*, 13(9).
- 47. Priyanka, J. H., & Parveen, N. (2024). DeepSkillNER: an automatic screening and ranking of resumes using hybrid deep learning and enhanced spectral clustering approach. *Multimedia Tools and Applications*, 83(16), 47503-47530.
- 48. Sathish, S., Thangavel, K., & Boopathi, S. (2010). Performance analysis of DSR, AODV, FSR and ZRP routing protocols in MANET. *MES Journal of Technology and Management*, 57-61.
- 49. Siva Prasad, B. V. V., Mandapati, S., Kumar Ramasamy, L., Boddu, R., Reddy, P., & Suresh Kumar, B. (2023). Ensemble-based cryptography for soldiers' health monitoring using mobile ad hoc networks. *Automatika: časopis za automatiku, mjerenje, elektroniku, računarstvo i komunikacije*, 64(3), 658-671.
- 50. Elechi, P., & Onu, K. E. (2022). Unmanned Aerial Vehicle Cellular Communication Operating in Nonterrestrial Networks. In *Unmanned Aerial Vehicle Cellular Communications* (pp. 225-251). Cham: Springer International Publishing.
- 51. Prasad, B. V. V. S., Mandapati, S., Haritha, B., & Begum, M. J. (2020, August). Enhanced Security for the authentication of Digital Signature from the key generated by the CSTRNG method. In 2020 Third International Conference on Smart Systems and Inventive Technology (ICSSIT) (pp. 1088-1093). IEEE.
- 52. Mukiri, R. R., Kumar, B. S., & Prasad, B. V. V. (2019, February). Effective Data Collaborative Strain Using RecTree Algorithm. In *Proceedings of International Conference on Sustainable Computing in Science, Technology and Management (SUSCOM), Amity University Rajasthan, Jaipur-India.*
- 53. Balaraju, J., Raj, M. G., & Murthy, C. S. (2019). Fuzzy-FMEA risk evaluation approach for LHD machine–A case study. *Journal of Sustainable Mining*, 18(4), 257-268.
- 54. Thirumoorthi, P., Deepika, S., & Yadaiah, N. (2014, March). Solar energy based dynamic sag compensator. In 2014 International Conference on Green Computing Communication and Electrical Engineering (ICGCCEE) (pp. 1-6). IEEE.
- 55. Vinayasree, P., & Reddy, A. M. (2025). A Reliable and Secure Permissioned Blockchain-Assisted Data Transfer Mechanism in Healthcare-Based Cyber-Physical Systems. *Concurrency and Computation: Practice and Experience*, 37(3), e8378.
- 56. Acharjee, P. B., Kumar, M., Krishna, G., Raminenei, K., Ibrahim, R. K., & Alazzam, M. B. (2023, May). Securing International Law Against Cyber Attacks through Blockchain Integration. In 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE) (pp. 2676-2681). IEEE.
- 57. Ramineni, K., Reddy, L. K. K., Ramana, T. V., & Rajesh, V. (2023, July). Classification of Skin Cancer Using Integrated Methodology. In *International Conference on Data Science and Applications* (pp. 105-118). Singapore: Springer Nature Singapore.
- 58. LAASSIRI, J., EL HAJJI, S. A. Î. D., BOUHDADI, M., AOUDE, M. A., JAGADISH, H. P., LOHIT, M. K., ... & KHOLLADI, M. (2010). Specifying Behavioral Concepts by engineering language of RM- ODP. *Journal of Theoretical and Applied Information Technology*, 15(1).
- 59. Prasad, D. V. R., & Mohanji, Y. K. V. (2021). FACE RECOGNITION-BASED LECTURE ATTENDANCE SYSTEM: A SURVEY PAPER. *Elementary Education Online*, 20(4), 1245-1245.
- 60. Dasu, V. R. P., & Gujjari, B. (2015). Technology-Enhanced Learning Through ICT Tools Using Aakash Tablet. In *Proceedings of the International Conference on Transformations in Engineering Education: ICTIEE* 2014 (pp. 203-216). Springer India.
- 61. Reddy, A. M., Reddy, K. S., Jayaram, M., Venkata Maha Lakshmi, N., Aluvalu, R., Mahesh, T. R., ... & Stalin Alex, D. (2022). An efficient multilevel thresholding scheme for heart image segmentation using a hybrid generalized adversarial network. *Journal of Sensors*, 2022(1), 4093658.
- 62. Srinivasa Reddy, K., Suneela, B., Inthiyaz, S., Hasane Ahammad, S., Kumar, G. N. S., & Mallikarjuna Reddy, A. (2019). Texture filtration module under stabilization via random forest optimization

- methodology. International Journal of Advanced Trends in Computer Science and Engineering, 8(3), 458-469.
- 63. Ramakrishna, C., Kumar, G. K., Reddy, A. M., & Ravi, P. (2018). A Survey on various IoT Attacks and its Countermeasures. *International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)*, 5(4), 143-150.
- 64. Sirisha, G., & Reddy, A. M. (2018, September). Smart healthcare analysis and therapy for voice disorder using cloud and edge computing. In 2018 4th international conference on applied and theoretical computing and communication technology (iCATccT) (pp. 103-106). IEEE.
- 65. Reddy, A. M., Yarlagadda, S., & Akkinen, H. (2021). An extensive analytical approach on human resources using random forest algorithm. *arXiv* preprint arXiv:2105.07855.
- 66. Kumar, G. N., Bhavanam, S. N., & Midasala, V. (2014). Image Hiding in a Video-based on DWT & LSB Algorithm. In *ICPVS Conference*.
- 67. Naveen Kumar, G. S., & Reddy, V. S. K. (2022). High performance algorithm for content-based video retrieval using multiple features. In *Intelligent Systems and Sustainable Computing: Proceedings of ICISSC 2021* (pp. 637-646). Singapore: Springer Nature Singapore.
- 68. Reddy, P. S., Kumar, G. N., Ritish, B., SaiSwetha, C., & Abhilash, K. B. (2013). Intelligent parking space detection system based on image segmentation. *Int J Sci Res Dev*, *1*(6), 1310-1312.
- 69. Naveen Kumar, G. S., Reddy, V. S. K., & Kumar, S. S. (2018). High-performance video retrieval based on spatio-temporal features. *Microelectronics, Electromagnetics and Telecommunications*, 433-441.
- Kumar, G. N., & Reddy, M. A. BWT & LSB algorithm based hiding an image into a video. *IJESAT*, 170-174.
- 71. Lopez, S., Sarada, V., Praveen, R. V. S., Pandey, A., Khuntia, M., & Haralayya, D. B. (2024). Artificial intelligence challenges and role for sustainable education in india: Problems and prospects. Sandeep Lopez, Vani Sarada, RVS Praveen, Anita Pandey, Monalisa Khuntia, Bhadrappa Haralayya (2024) Artificial Intelligence Challenges and Role for Sustainable Education in India: Problems and Prospects. Library Progress International, 44(3), 18261-18271.
- 72. Yamuna, V., Praveen, R. V. S., Sathya, R., Dhivva, M., Lidiya, R., & Sowmiya, P. (2024, October). Integrating AI for Improved Brain Tumor Detection and Classification. In 2024 4th International Conference on Sustainable Expert Systems (ICSES) (pp. 1603-1609). IEEE.
- 73. Kumar, N., Kurkute, S. L., Kalpana, V., Karuppannan, A., Praveen, R. V. S., & Mishra, S. (2024, August). Modelling and Evaluation of Li-ion Battery Performance Based on the Electric Vehicle Tiled Tests using Kalman Filter-GBDT Approach. In 2024 International Conference on Intelligent Algorithms for Computational Intelligence Systems (IACIS) (pp. 1-6). IEEE.
- 74. Sharma, S., Vij, S., Praveen, R. V. S., Srinivasan, S., Yadav, D. K., & VS, R. K. (2024, October). Stress Prediction in Higher Education Students Using Psychometric Assessments and AOA-CNN-XGBoost Models. In 2024 4th International Conference on Sustainable Expert Systems (ICSES) (pp. 1631-1636). IEEE.
- 75. Anuprathibha, T., Praveen, R. V. S., Sukumar, P., Suganthi, G., & Ravichandran, T. (2024, October). Enhancing Fake Review Detection: A Hierarchical Graph Attention Network Approach Using Text and Ratings. In 2024 Global Conference on Communications and Information Technologies (GCCIT) (pp. 1-5). IEEE.
- 76. Shinkar, A. R., Joshi, D., Praveen, R. V. S., Rajesh, Y., & Singh, D. (2024, December). Intelligent solar energy harvesting and management in IoT nodes using deep self-organizing maps. In 2024 International Conference on Emerging Research in Computational Science (ICERCS) (pp. 1-6). IEEE.
- 77. Praveen, R. V. S., Hemavathi, U., Sathya, R., Siddiq, A. A., Sanjay, M. G., & Gowdish, S. (2024, October). AI Powered Plant Identification and Plant Disease Classification System. In 2024 4th International Conference on Sustainable Expert Systems (ICSES) (pp. 1610-1616). IEEE.
- 78. Dhivya, R., Sagili, S. R., Praveen, R. V. S., VamsiLala, P. N. V., Sangeetha, A., & Suchithra, B. (2024, December). Predictive Modelling of Osteoporosis using Machine Learning Algorithms. In 2024 4th International Conference on Ubiquitous Computing and Intelligent Information Systems (ICUIS) (pp. 997-1002). IEEE.
- 79. Kemmannu, P. K., Praveen, R. V. S., Saravanan, B., Amshavalli, M., & Banupriya, V. (2024, December). Enhancing Sustainable Agriculture Through Smart Architecture: An Adaptive Neuro-Fuzzy Inference System with XGBoost Model. In 2024 International Conference on Sustainable Communication Networks and Application (ICSCNA) (pp. 724-730). IEEE.
- 80. Praveen, R. V. S. (2024). Data Engineering for Modern Applications. Addition Publishing House