

# A Hybrid Deep Neural Network for Multimodal Deepfake Detection

Cinana Vinod<sup>1</sup>, Fathimathu Nasna SP<sup>2</sup> JishnuPrasadKV<sup>3</sup>, Shifad Shukkoor<sup>4</sup>,Ms.Vanimol Sajan<sup>5</sup> Dept. of Computer Science and Engineering

Vimal Jyothi Engineering College, Chemperi, Kannur, India *cinanavinod8@gmail.com*<sup>1</sup>, *fathimathunasna4540@gmail.com*<sup>2</sup>, *jpkv2002@gmail.com*<sup>3</sup>, *shifadshukkoor123@gmail.com*<sup>4</sup>, *vanimolsajan@vjec.ac.in*<sup>5</sup>

*Abstract*—Detecting AI-manipulated media, especially deep- fake images, videos, and audio, is critical in combating privacy and security threats. This project introduces a hybrid deep learning model that integrates content and trace feature extrac- tion to enhance detection accuracy and robustness. By utilizing CNN and RNN architectures, the model processes datasets from DeepFake and Face2Face algorithms, effectively identifying subtle manipulations even under challenging conditions like video compression. Experimental results highlight its superior perfor- mance over baseline methods, supported by Class Activation Maps (CAMs) that pinpoint critical areas.Additionally, this model integrates seamlessly with a social media-like application to provide real-time monitoring of uploaded media. The system not only detects and flags AI-generated manipulations but also alerts users about suspicious content as it is shared or uploaded. This real-time functionality enhances the system's usability in dynamic environments, making it an efficient tool for advancing media forensics, combating misinformation, and ensuring content authenticity across digital platforms.

Keywords: Deepfake Detection, Hybrid Neural Net- works, Multimodal Analysis, CNN, RNN, Media Forensics, Face2Face.

## I. INTRODUCTION

Deepfake technology, which uses AI to generate highly realistic synthetic media, presents significant challenges in the digital era. Techniques like DeepFake and Face2Face enable the sophisticated manipulation of images, videos, and audio, often without the knowledge or consent of the individuals portrayed. This raises serious concerns about misinformation, impersonation, and reputation damage, as such manipulations can be used maliciously to deceive the public or harm individuals. The rapid advancements in AI-driven content generation have far outpaced traditional detection methods, making it increasingly difficult to distinguish between real and manipulated media. Conventional detection techniques often struggle with accuracy, especially when dealing with compressed video formats or subtle alterations that evade basic forensic analysis, contributing to a growing mistrust in the authenticity of digital content. This project seeks to address these challenges by leveraging state-of-the-art deep learning techniques to develop a hybrid neural network capable of accurately detecting AI-generated manipulations. The system integrates multimodal analysis, processing images, videos, and audio simultaneously for a more comprehensive and reliable deepfake detection solution. Unlike traditional methods thatfocus on single modalities, this approach enhances detection accuracy by identifying manipulation traces across various types of media. Convolutional Neural Networks (CNNs) are employed for image and video analysis, while Recurrent Neural Networks (RNNs) are utilized for audio detection. To further enhance interpretability, the system uses Class Activation Maps (CAMs) to visually highlight manipulated regions, making it easier for forensic analysts to verify content authenticity. The hybrid approach ensures high accuracy and resilience against evolving deepfake techniques. By incorpo- rating both content-based and forensic trace feature extraction methods, the system can detect visible and subtle artifacts introduced during the deepfake generation process. This guar- antees comprehensive media analysis, even under challenging conditions such as video compression, low resolution, or adversarial perturbations. Additionally, the system is designed to be scalable and userfriendly, providing intuitive visual feedback such as CAMs for images and spectrograms for audio. This makes it accessible not only to forensic experts but also to general users who may need to verify media authenticity. The system's modular design also allows for easy integration with existing forensic software and adapta- tion to new deepfake manipulation techniques as technology continues to advance. To further enhance accessibility, this project envisions integrating the detection framework into a mobile app. This app will allow users to upload images, videos, and audio directly from their devices for real-time deepfake detection. It will provide instant feedback through visual



cues like CAMs for images and spectrograms for audio, making it a valuable tool for individuals, journalists, or organizations seeking to verify digital content on the go. The app will also support integration with existing forensic tools and digital security systems, ensuring that it aligns with current efforts to combat misinformation. Beyond improving detection capabilities, this project significantly contributes to the field of media forensics by providing a practical tool for combating the spread of misinformation. By enhancing digital content verification, the system supports efforts to maintain trust in online media, helping platforms, journalists, and law enforcement agencies identify and mitigate the impact of deepfake-based deception. Furthermore, the system aligns with ethical and legal frameworks advocating for responsible AI use. It assists regulatory bodies and organizations in enforcing

policies aimed at curbing the spread of manipulated content, ensuring that digital media remains a trusted source of infor- mation. As deepfake technology evolves, the proposed system offers a proactive solution to safeguard digital integrity and ensure that the detection of synthetic media keeps pace with rapid advancements in AI-driven content creation.

## II. RELATED WORKS

FDT: A Python Toolkit for Fake Image and Video Detec- tion introduces the Fake Detection Toolkit (FDT), a Pythonbased toolkit designed to identify and visualize manipulated images and videos. The toolkit provides a streamlined detection process, making it accessible to a broader audience, including users on social media platforms like Twitter.FDT combines both content-based and forensic feature extraction techniques to identify inconsistencies in image and video data. The toolkit's integration with social media is particularly valuable for addressing misinformation in real-time.FDT offers an intuitive interface that allows non-experts to detect and visualize manipulated media, empowering users with the tools to counter fake content independently. This paper highlights the growing need for accessible deepfake detection technologies that can be easily integrated with social media platforms and addresses the challenge of widespread misin- formation by providing a user-friendly solution for identifying and analyzing manipulated media. [1] Deepfake Detection in Digital Media Forensics study presents a deepfake detection model that leverages ResNext and Long Short-Term Memory (LSTM) networks to capture spatial and temporal inconsis- tencies in video content. By combining these two types of neural networks, the model is particularly effective at handling sequential video data. The ResNext component processes individual frames for spatial anomalies, while the LSTM layer analyzes temporal inconsistencies, enabling the model to identify subtle manipulation artifacts over time. Tested on the Celeb-DF dataset, this model achieves a high detection accuracy of 91%, demonstrating its effectiveness in real-world applications. The study also discusses the challenges posed by video compression, which often reduces the visibility of tampering artifacts. The proposed approach, however, maintains strong detection performance even under compressed conditions, underscoring its robustness. This paper highlights the importance of handling both spatial and temporal features in video-based deepfake detection, as well as the value of using multimodal neural networks to address the unique challenges posed by manipulated video content. [2]

JRC: Deepfake Detection via Joint Reconstruction and Classification introduces the Joint Reconstruction and Classification (JRC) model is introduced, which uniquely combines reconstruction and classification tasks in a single framework to enhance feature learning for deepfake detection. The JRC model achieves state-of-the-art performance across multiple deepfake datasets, demonstrating its ability to detect both obvious and subtle manipulations. By reconstructing manipulated faces, the model is able to identify anomalies in the reconstruction process, providing an additional layerof evidence for classification decisions. The classification component further refines these decisions by identifying inconsistencies in the reconstructed facial features. This dual approach enables the model to outperform traditional detection techniques, particularly when handling diverse datasets or media with significant compression. JRC's robustness in the face of video compression and subtle manipulations makes it a valuable tool for practical applications where manipulated media often undergoes additional processing before distribution. [3] Exposing Deepfakes Using a Deep Multilayer Perceptron–Convolutional Neural Network Model introduce a hybrid model that combines a multilayer perceptron (MLP) with a CNN to enhance deepfake detection capabilities. The MLP effectively handles high-dimensional data, while the CNN specializes in extracting spatial features from images, making this combination particularly effective for identifying manipulated facial characteristics. The model reaches an accuracy of 84% and an area under the curve (AUC) score of 0.87, demonstrating a



balanced and effective performance across a range of datasets. The hybrid approach allows the model to capture a diverse set of features, enabling it to detect manipulation across various types of media and under different levels of video compression. The study suggests that this hybrid MLP- CNN architecture provides a more comprehensive detection approach, capturing both detailed and abstract patterns in manipulated media that single-architecture models often miss. This model is particularly well-suited for applications requiring high accuracy across diverse types of fake media.

[4] DLFMNet: End-to-End Detection and Localization of Face Manipulation Using Multi-Domain Features by Chen et al.(2021) presents a novel approach to detecting and localizing face manipulations, such as those created by deepfake techniques. The proposed method, DLFMNet, is an end-to-end deep learning model that utilizes multi-domain features to improve detection performance. Rather than relying on a single type of feature, DLFMNet incorporates a variety of face-related features, such as texture, geometry, and lighting information, to robustly detect and pinpoint manipulations in facial images and videos. DLFMNet's architecture is designed to handle both the detection and localization of manipulated regions within the face, making it capable of not only identifying whether an image has been altered but also precisely locating the areas affected by manipulation. The authors demonstrate that their method significantly outperforms existing face manipulation detection models, especially in terms of accuracy and robustness, by leveraging these multi-domain features. Through experiments conducted on several benchmark datasets, the model shows strong performance in real-world scenarios, and the paper concludes that DLFMNet offers a promising solution for face manipulation detection in multimedia content. [8]

Ref No.	Paper Title	Author(s)	Methods	Advantages	Disadvantages		
1	FDT: A	Sharma et al.	Fake Detection Toolkit	Easy-to-use, real-time	Limited		
	Python		(FDT): identifies and	integration, suitable for	datase		
	Toolkit for		visualizes manipulated	non-experts	t coverage, less		
	Fake Image		media; integrates with		effective against		
	and		social media for real-		sophisticated		
	Vide		time misinformation		manipulations		
	o Detection		detection.				
2	Deepfake	Singh et al.	Combines ResNext and	High accuracy, robust	Requires large		
	Detection in		LSTM networks to de-	under	computational		
	Digital Media		tect spatial and tem-	compression	resources,		
	Forensics		poral inconsistencies in	,	drop		
			deepfake videos; 91%	handles	s with unseen styles		
			detection accuracy on	tempora			
			Celeb-DF.	l inconsistencies			
3	JRC:	Zhang et al.	JRC model enhances	Handles diverse	High computational		
	Deepfake		feature learning via re-	datasets,	cost,		
	Detection via		construction and clas-	effectiv	comple		
	Joint		sification; excels with	e under compression,	x architecture		
	Recon		compressed and	combines			
	- struction		diverse datasets.	feature			
	and			learning with			
	Classification			reconstruction			
4	Exposing	Kumar et al.	Combines MLP and	Balanced performance,	Lower accuracy than		
	Deepfakes		CNN for deepfake de-	hybrid	advanced models, lim-		
	Using		tection; achieves 84%	capture	ited robustness under		
	MLP		accuracy and 0.87	s abstract and spatial	severe compression		
	- CNN		AUC across diverse	features			
			datasets.				



International Journal of Scientific Research in Engineering and Management (IJSREM)

Volume: 09 Issue: 04 | April - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

5	DLFMNet:	Chen et al.	DLFMNet uses 1	multi-	Detects	and	localizes	High	comp	utational
	End-to-End		domain features	s to	manipula	ations,	, robust	cost,	may	require
	Detection and		detect and lo	calize	in real-	-world	l cases,	specialized hardware		
	Localization		manipulated	face	uses	mult	i-domain			
			regions; robust	in	features					
			real-world scenarios.							

### III. PROPOSED METHODOLOGY

The proposed Hybrid Deepfake Detection System is an advanced AI-driven framework designed to detect manipulated media with high accuracy by leveraging multimodal anal- ysis techniques. This system integrates sophisticated neural network architectures to analyze and identify AI-generated fake content, such as images, videos, and audio, created using techniques like DeepFake and Face2Face. The system architecture comprises four main modules: the Image Analysis Module, Video Analysis Module, Audio Analysis Module, and Multimodal Fusion Module.

The Image Analysis Module employs Convolutional Neural Network (CNN) architectures, such as ResNet-18, to extract fine-grained facial features, identifying anomalies like unnatural textures and lighting inconsistencies within static images. For video analysis, the Video Analysis Module utilizes a combination of CNNs and Recurrent Neural Networks (RNNs) or 3D CNNs to process sequential video

frames. This module tracks facial movements and evaluates temporal consistency to detect frame-by-frame manipulations, including subtle identity swaps or expression changes. The Audio Analysis Module processes speech patterns using CNNs and RNNs, relying on Mel-frequency cepstral coefficients (MFCCs) and other features to detect irregularities such as robotic tones or abrupt tonal shifts in synthesized audio. The Multimodal Fusion Module integrates features extracted from the image, video, and audio analysis modules through a fully connected neural network, enabling a comprehensive and robust authenticity assessment across diverse media types.

The system is designed with a layered architecture to ensure seamless interaction among its components and scalability for real-world applications. The frontend design includes a user-friendly interface allowing general users to upload media for verification, an analyst dashboard providing advanced visualization and comparison tools, and an admin interface for managing system configurations. The backend architec- ture, implemented using Python frameworks and a MySQL database, supports efficient data processing, storage, and re-

trieval. It incorporates asynchronous pipelines to handle high- volume media processing and optimized data handling for real- time detection.

At the core of the system are specialized deep learning models tailored for each media type. These include CNN classifiers for detecting manipulated facial features in images, 3D CNNs with RNN layers for evaluating temporal inconsis- tencies in video frames, and RNNs for analyzing sequential patterns in audio. The multimodal fusion network consolidates the extracted features, enhancing detection accuracy by cross- referencing information from all three modalities.

To ensure reliability and transparency, the system integrates evaluation and visualization tools. Performance metrics, in- cluding accuracy, F1-score, and Area Under Curve (AUC), are employed to assess and refine model performance. Class Acti- vation Maps (CAMs) are used to visualize regions of focus in images and video frames, providing insights into the model's decision-making process. Furthermore, the system undergoes robustness testing against media compression, particularly with the H.264 codec, to ensure consistent performance in typical internet scenarios. This comprehensive methodology addresses the growing challenges posed by AI-generated ma- nipulated media, providing an effective and scalable solution for deepfake detection.



Volume: 09 Issue: 04 | April - 2025

SJIF Rating: 8.586



Fig. 1. Proposed architecture diagram of hybrid deep neural network model

The Fig.1 offers an overview of the Hybrid Deepfake Detection System, showcasing its components, interactions, and data flow for image, video, and audio analysis. It high- lights the deployment structure, interaction pathways, and data processing pipeline, providing a clear understanding of the system's functionality. Here is a breakdown of the system architecture as described:

## A. User Roles and Functions

• General Users:Primarily access media verification func- tionalities, where they can upload images, audio, or video files to receive deepfake authenticity assessments. Users benefit from real-time feedback, authenticity scores, and visual indicators that highlight manipulated areas.

- Forensic Analysts:For professionals conducting in-depth media analysis, this role provides advanced features, including access to Class Activation Maps (CAMs) for examining manipulation hot spots and detailed metrics on detection accuracy and performance. Administrators:Responsible for overseeing system oper- ation, ensuring data accuracy, and managing user ac- counts. Administrators can modify system parameters, adjust detection thresholds, and monitor overall system performance to ensure optimal functionality.

# B. Detection and Analysis Modules

• Image Detection (Convolutional Neural Network): This module leverages a CNN (e.g., ResNet-18) to analyze static images. By extracting critical facial features, it identifies inconsistencies in color, texture, and structure associated with image-based deepfakes. The CNN pro- cesses aligned face data, focusing on identifying subtle artifacts.

• Video Detection (3D CNN and RNN): This module pro- cesses sequential video frames to assess facial consistency across time. A 3D CNN captures spatiotemporal features, while an RNN tracks continuity in facial expressions and movements. This dual approach identifies irregularities in manipulated videos, particularly effective for detecting deepfake-generated expressions and identity swaps.

• Audio Detection (RNN for Temporal Analysis):Focuses on analyzing voice patterns using an RNN. The model evaluates audio features, such as MFCCs, to identify synthetic speech patterns and unnatural tonal shifts. This module aids in identifying audio-based deepfake content commonly found in political and celebrity impersonation videos.

• Multimodal Fusion Network: This integrative module combines features extracted from image, video, and audio models, enhancing overall detection accuracy by analyz- ing multiple modalities concurrently. Using a fully connected neural network, it merges results from individual models, providing a holistic authenticity score based on cross-referenced data.

# C. Preprocessing Components

Preprocessing in the Hybrid Deepfake Detection System involves preparing image, video, and audio data for accurate analysis. For images, this includes facial detection and align- ment, resizing, and color normalization to ensure uniformity, enabling the CNN to focus on key facial regions and improving manipulation detection. For videos, key frames are extracted at intervals, converted into a sequence of images, and aligned to allow the 3D CNN to assess spatiotemporal changes accurately within facial regions. Audio data is segmented into short, consistent clips to facilitate feature extraction, with Mel- frequency cepstral coefficients (MFCCs) and other features computed to help



the RNN model identify speech irregularities typical of synthesized audio content.

## D. Output Interfaces

The detection output interface presents the results of the detection models, providing authenticity scores and visual highlights, such as Class Activation Maps (CAMs) for images

and videos. These outputs enable users to quickly identify manipulated areas and assess confidence levels for authenticity evaluations. Additionally, the system generates detailed analysis reports that include comprehensive detection data, such as model metrics, feature importance, and manipulation locations, which can be downloaded for deeper analysis or use as evidence in digital forensics. The system also incorporates user feedback into a learning loop, allowing for adjustments to detection thresholds based on real-world performance and user inputs, thereby enhancing adaptability to evolving deepfake techniques.



## Fig. 2. Dataflow diagram

The Fig.2 outlines the process of detecting manipulated media, starting with user input, where images, videos, or audio are uploaded through an interface. The media undergoes pre- processing, including face detection and alignment for images and videos, frame extraction for videos, and segmentation with MFCC feature extraction for audio. Features are then extracted using CNNs for visual data and neural networks for audio. The hybrid model performs multimodal fusion to combine features from all media types, which are analyzed to classify the content as "real" or "fake." Results are displayed with authenticity scores, CAMs for visualizing manipulated regions, and spectrograms for audio, ensuring transparency in the detection process.

#### IV. EXPERIMENTAL RESULTS

We first provide dataset description. Then, we present the baselines, the training details, and evaluation metrics.

#### A. Datasets

For the experiments, we used three fake face datasets that were created by using Face2Face ,FakeAudio dataset and DeepFake techniques. The details for each dataset are as follows.

1) *Face2Face Dataset:* The dataset includes video sam- ples manipulated using expression-based techniques like Face2Face. The dataset contains videos sourced at a resolution of 480p or higher, ensuring good quality for analysis. Frames are extracted from these videos to serve as training samplesfor the detection model. The Face2Face dataset is specifically valuable for detecting expression manipulations, where facial movements and expressions are altered while retaining the original identity.

2) *DeepFake Dataset:* The dataset was created using an autoencoder-based DeepFake generation technique, involving identity-swapping manipulations. Videos of multiple individ- uals were used, ensuring diverse face shapes, genders, and attributes to train the model effectively. Frames were extracted from these manipulated videos, and facial alignment was applied to maintain consistency. The dataset also includes audio clips with AI-synthesized



voices, which mimic natural speech patterns while introducing subtle irregularities in tone and pitch. This custom dataset provides a rich variety of scenarios for evaluating the model's robustness across different manipulation techniques.

*3) Fake Audio datasets:* To address audio-based deepfake detection, the dataset includes AI-synthesized audio clips that exhibit common manipulation artifacts such as unnatural tonal shifts or robotic sounds. The dataset is processed to extract features like Mel-frequency cepstral coefficients (MFCCs), which help capture the unique characteristics of synthesized speech. These clips enable the model to differentiate between genuine and manipulated voice patterns effectively.

These datasets collectively cover a wide range of manipulation techniques for images, videos, and audio, allowing the hybrid model to generalize across diverse scenarios and media types.

#### B. Baseline Models

1) ResNet-18: ResNet-18 is a convolutional neural network (CNN) widely used for image and video analysis tasks. It is particularly effective in extracting high-level content-based features, such as facial textures, lighting, and structural in- consistencies, which are crucial for detecting manipulations. Its residual connections help prevent the vanishing gradient problem, making it efficient for deepfake detection tasks. In this project, ResNet-18 serves as a strong benchmark for evaluating the effectiveness of the proposed hybrid system.

2) *MesoNet:* MesoNet is a shallow CNN model designed specifically for detecting facial manipulations in images and videos. Its architecture is optimized for identifying subtle pixel-level artifacts, such as unnatural textures or blend- ing inconsistencies, introduced during deepfake generation. MesoNet's lightweight design ensures computational effi- ciency, making it suitable for applications requiring quick and reliable manipulation detection.

*3) VGG-16:* VGG-16, a deep CNN architecture with 16 layers, is used to extract spatial features from manipulated media. It identifies patterns such as texture anomalies, lighting irregularities, and facial distortions, which are common in deepfake content. Despite its effectiveness in handling basic manipulations, VGG-16 may struggle with advanced manipu- lation techniques and is computationally heavier compared to more modern architectures.

4) Long Short-Term Memory (LSTM): LSTM networks are a type of recurrent neural network (RNN) designed to process sequential data, such as video frames or audio recordings. In this project, LSTMs are employed to capture temporal dependencies and detect inconsistencies in frame sequences or speech patterns. By analyzing subtle changes over time, LSTMs are highly effective for identifying deepfake manipu- lations in videos and audio.

5) Spatial Rich Model (SRM) with SVM: The Spatial Rich Model (SRM) is a traditional image forensics technique that extracts handcrafted features to detect manipulation artifacts, such as noise inconsistencies or resampling traces. These features are classified using a Support Vector Machine (SVM) to categorize media as real or fake. While effective for basic forgery detection, SRM+SVM lacks the flexibility of deep learning models and struggles to handle high-quality manipulations.

#### C. Training Details

The training process for the hybrid deepfake detection system involved preprocessing media inputs, including facial alignment and resizing for images and videos, and MFCC fea- ture extraction for audio. The architecture consisted of ResNet- 18 for images, 3D CNNs and LSTMs for video analysis, and a CNN-RNN model for audio, with all features integrated through a multimodal fusion network. Training was conducted using TensorFlow and PyTorch with a learning rate of 0.001, batch size of 32, and the Adam optimizer over 50 epochs, in- corporating regularization techniques like dropout, early stop- ping, and data augmentation. The models were evaluated using accuracy, F1-score, and robustness tests under compression, and the process was executed on high-performance GPUs, ensuring efficient handling of large datasets and achieving a robust, multimodal detection framework.



## V. CONCLUSION AND FUTURE WORK

The Hybrid Deep Neural Network for Multimodal Deepfake Detection presents a cutting-edge solution to the growing threat of AI-manipulated media by combining advanced im- age, video, and audio analysis techniques. This system effec- tively addresses challenges like subtle manipulation artifacts, video compression, and diverse media formats, delivering robust detection capabilities. By integrating CNN, 3D CNN, and RNN models alongside multimodal fusion, the system ensures comprehensive and accurate identification of manipu- lated content. It also enhances user experience through intu- itive interfaces and visualization tools, making it practical for applications in media forensics, security, and misinformation prevention.

This project underscores the potential of leveraging AI- driven solutions to tackle emerging challenges in digital media authentication. By automating detection processes and offering high reliability, it establishes a strong foundation for combat- ing the misuse of synthetic media. Future work will focus on expanding the system's capabil- ities to handle real-time detection in live streaming scenarios and adapting it for multilingual and cross-cultural media. Enhancing robustness against evolving deepfake technologies, such as generative adversarial networks (GANs), will also be a priority. Additionally, integrating this system into larger frameworks for social media platforms and digital forensics can further streamline efforts to combat misinformation and improve global media security. This progressive approach ensures the system remains adaptive and scalable, paving the way for its continued relevance in the fight against AI- manipulated content.

#### REFERENCES

[1] Raj, Surbhi and Mathew, Jimson and Mondal, Arijit, "FDT: A Python Toolkit for Fake Image and Video Detection," *SoftwareX*, vol. 22, pp. 101395, 2023, Elsevier.

<sup>[2]</sup> Vamsi, Vurimi Veera Venkata Naga Sai and Shet, Sukanya S and Reddy, Sodum Sai Mohan and Rose, Sharon S and Shetty, Sona R and Sathvika, S and Supriya, MS and Shankar, Sahana P, "Deepfake Detection in Digital Media Forensics," *Global Transitions Proceedings*, vol. 3, no. 1,

pp. 74-79, 2022, Elsevier.

[3] Yan, Bosheng and Li, Chang-Tsun and Lu, Xuequan, "JRC: Deepfake Detection via Joint Reconstruction and Classification," *Neurocomputing*,

pp. 127862, 2024, Elsevier.

[4] Kolagati, Santosh and Priyadharshini, Thenuga and Rajam, V Mary Anita, "Exposing Deepfakes Using a Deep Multilayer Perceptron– Convolutional Neural Network Model," *International Journal of Infor- mation Management Data Insights*, vol. 2, no. 1, pp. 100054, 2022, Elsevier.

<sup>[5]</sup> Camara, Mateus Karvat and Postal, Adriana and Maul, Tomas Henrique and Paetzold, Gustavo Henrique, "Can Lies Be Faked? Comparing Low- Stakes and High-Stakes Deception Video Datasets from a Machine Learning Perspective," *Expert Systems with Applications*, vol. 249, pp. 123684, 2024, Elsevier.

<sup>[6]</sup> Guo, Zhiqing and Yang, Gaobo and Chen, Jiyou and Sun, Xingming, "Fake Face Detection via Adaptive Manipulation Traces Extraction Network," *Computer Vision and Image Understanding*, vol. 204, pp. 103170, 2021, Elsevier.

[7] Ju, Xingwang, "An Overview of Face Manipulation Detection," *Journal of Cybersecurity*, vol. 2, no. 4, pp. 197, 2020, Tech Science Press.

<sup>[8]</sup> Chen, Peng and Liu, Jin and Liang, Tao and Yu, Cai and Zou, Shuqiao and Dai, Jiao and Han, Jizhong, "DLFMNet: End-to-End Detection and Localization of Face Manipulation Using Multi-Domain Features," 2021 IEEE International Conference on Multimedia and Expo (ICME), pp. 1– 6, 2021, IEEE.

[9] Fernando, Tharindu and Fookes, Clinton and Denman, Simon and Sridharan, Sridha, "Detection of Fake and Fraudulent Faces via Neural Memory Networks," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 1973–1988, 2020, IEEE.

[10] Bekci, Burak and Akhtar, Zahid and Ekenel, Hazım Kemal, "Cross- Dataset Face Manipulation Detection," 2020 28th Signal Processing and Communications Applications Conference (SIU), pp. 1–4, 2020, IEEE.



[11] Chang, Xu and Wu, Jian and Yang, Tongfeng and Feng, Guorui, "Deep- Fake Face Image Detection Based on Improved VGG Convolutional Neural Network," *2020 39th Chinese Control Conference (CCC)*, pp. 7252–7256, 2020, IEEE.

[12] Li, Gen and Cao, Yun and Zhao, Xianfeng, "Exploiting Facial Symmetry to Expose Deepfakes," 2021 IEEE International Conference on Image Processing (ICIP), pp. 3587–3591, 2021, IEEE.

<sup>[13]</sup> Wang, Xinyao and Yao, Taiping and Ding, Shouhong and Ma, Lizhuang, "Face Manipulation Detection via Auxiliary Supervision," *Neural In- formation Processing: 27th International Conference, ICONIP 2020, Bangkok, Thailand, November 23–27, 2020, Proceedings, Part I 27*, pp. 313–324, 2020, Springer.

[14] Lai, Yingxin and Yang, Guoqing and He, Yifan and Luo, Zhiming and Li, Shaozi, "Selective Domain-Invariant Feature for Generalizable Deep- fake Detection," *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2335–2339, 2024, IEEE.