



International Journal of Engineering Research and Science & Technology

www.ijerst.org

ISSN : 2319-5991

Vol. 22 No. 2 (2026)



ijerst.editor@gmail.com
editor@ijerst.com

Research Paper

DEEFAKE IMAGE RECOGNITION VIA FEATURE LEARNING IN CONVOLUTIONAL NEURAL NETWORKS

Avanapu Aparna, Chadaram Sri Harsha Vardhan,

Boni Appanna Nivas, Ellabelli Priyatham

Department of Computer Science & Engineering (DS)

Raghu Engineering College

Visakhapatnam, Andhra Pradesh, India

avanapuaparna8@gmail.com

Mr. V. Vidya Sagar

Asst. Prof., Department of Computer Science & Engineering (DS)

Raghu Engineering College

Visakhapatnam, Andhra Pradesh, India

sagarpoornavivek@gmail.com

Abstract—The surge in synthetically generated facial imagery—commonly termed deepfakes—poses escalating threats to information integrity and digital forensics. This paper presents a Convolutional Neural Network (CNN) framework employing transfer learning from the MobileNet backbone, pre-trained on ImageNet, for binary classification of real versus GAN-synthesized facial images. A lightweight classification head comprising Global Average Pooling, a Dense(128) layer with ReLU activation, Dropout(0.5) regularization, and a Softmax output is fine-tuned on the Kaggle Deepfake-and-Real-Images benchmark dataset. Training for five epochs with the Adam optimizer ($\text{lr} = 3 \times 10^{-4}$) achieves 96.1% validation accuracy and a weighted F1-score of 0.960, surpassing several published baselines. An AUC-ROC of 0.991 confirms strong discriminative calibration. A Flask-based REST API and React.js frontend are deployed to enable real-time, browser-accessible inference.

Keywords—Deepfake Detection; Convolutional Neural Networks; Transfer Learning; MobileNet; GAN; Image Forensics; Flask; React.js

I. INTRODUCTION

The rapid proliferation of synthetically generated images—commonly referred to as deepfakes—poses serious threats to information integrity, digital forensics, and public trust. Advances in Generative Adversarial Networks (GANs) and diffusion-based synthesis techniques have dramatically lowered the barrier for creating hyper-realistic forged facial imagery, enabling misuse in disinformation campaigns, identity fraud, and non-consensual media manipulation.

Traditional image forensics relied on pixel-level statistical anomalies, JPEG compression artifacts, and noise patterns to detect manipulation. However, these handcrafted feature approaches fail against modern neural synthesis pipelines that produce perceptually indistinguishable fakes. There is therefore a pressing need for robust, end-to-end learning-based detectors that automatically discover discriminative features at multiple semantic levels.

This paper presents a Convolutional Neural Network (CNN) framework for deepfake image classification that leverages transfer learning from the MobileNet backbone pre-trained on ImageNet. By freezing the convolutional feature extractor and fine-tuning a compact classification head, our approach achieves 96% validation accuracy on the Kaggle Deepfake-and-Real-Images benchmark dataset. A Flask-based REST API and React.js frontend are deployed to demonstrate real-time inference. The remainder of this paper is organized as follows: Section II reviews related literature; Section III describes the dataset; Section IV details the methodology; Section V presents experimental results; Section VI discusses explainability; Section VII describes the system architecture; Section VIII concludes.

II. RELATED WORK

Early deepfake detection methods exploited biological signals and physical inconsistencies. Li et al. [1] demonstrated that GAN-generated faces exhibit characteristic blending boundaries when the forged region is warped onto a target face. Similarly, Matern et al. [2] identified low-level visual artifacts—notably irregular eye reflections and irregular teeth textures—as reliable forgery indicators accessible without any learned model. These handcrafted approaches, while interpretable, showed severe performance degradation as synthesis quality improved.

Deep learning detectors emerged as synthesis techniques matured. Rossler et al. [3] introduced FaceForensics++, a large-scale benchmark containing over one million manipulated video frames from four manipulation methods (DeepFakes, Face2Face, FaceSwap, NeuralTextures). They showed that binary CNN classifiers trained on compressed frames achieve up to 99% accuracy at low compression levels, but accuracy degrades significantly at high compression ratios encountered in social-media distribution. Nguyen et al. [4] proposed CapsuleForensics, replacing conventional fully connected layers with dynamic routing

capsule layers to improve robustness to unseen manipulation types.

Transfer learning has been widely adopted to mitigate the limited availability of labelled deepfake data. Dang et al. [5] fine-tuned Xception and ResNet backbones pre-trained on ImageNet, achieving over 97% accuracy on FaceForensics++ with minimal domain-specific training. Our work similarly exploits the representational power of pre-trained features but employs the lightweight MobileNet architecture, which reduces computational cost while maintaining competitive accuracy—a property desirable for web-deployable inference scenarios. Unlike prior works that target video frame sequences, we focus on single-image classification, evaluating the model on a balanced still-image dataset.

III. DATASET

A. Source and Composition

The experiments employ the Kaggle "Deepfake and Real Images" dataset, which contains a balanced collection of face images drawn from two classes: Real and Fake. Real images are sourced from the Flickr Faces High Quality (FFHQ) dataset, comprising photographs of real human subjects captured under diverse lighting conditions, poses, and ethnicities. Fake images are GAN-synthesized portraits generated by StyleGAN2, representing state-of-the-art synthesis quality as of the dataset release.

The dataset is partitioned into a training split located at Dataset/Train/ and a separate test split. For our experiments, the training split is further divided by the ImageDataGenerator into 75% training (approx. 14,000 images) and 25% validation (approx. 4,700 images) subsets. Each subset contains approximately equal numbers of Real and Fake images, ensuring balanced class distribution and avoiding accuracy-inflation bias from class imbalance.

B. Preprocessing Pipeline

All images are resized to 256×256 pixels using bilinear interpolation to match the MobileNet input requirement. Pixel values are rescaled to the $[0, 1]$ range via the $\text{rescale}=1./255$ factor within the Keras ImageDataGenerator. During training, mild augmentation is applied: random zoom up to 20% of the image width ($\text{zoom_range}=0.2$), with any newly introduced pixel regions filled using the nearest-neighbor strategy ($\text{fill_mode}=\text{'nearest'}$). No horizontal flipping or rotation is applied in order to preserve facial orientation cues exploited by the detector. Validation and test images undergo only rescaling without augmentation to ensure a fair evaluation.

TABLE I. DATA CONFIGURATION PARAMETERS

Parameter	Value
Input image size	$256 \times 256 \times 3$
Batch size	32
Training split	75% (~14,000 images)
Validation split	25% (~4,700 images)
Class distribution	Balanced (Real / Fake)
Augmentation	Zoom $\pm 20\%$, fill_mode=nearest

Pixel normalization	Rescale $1/255 \rightarrow [0, 1]$
---------------------	------------------------------------

C. Dataset Integrity

The Fake class exclusively contains StyleGAN2-generated images, which exhibit a distinct generative fingerprint different from autoencoder-based face-swap methods used in video deepfakes. This specialization means the model primarily learns features discriminating photographic facial texture from neural-synthesis artifacts, such as grid-like high-frequency patterns in the GAN output spectrum, unrealistic hair strand rendering, and inconsistent peripheral background blending.

IV. METHODOLOGY

A. Transfer Learning with MobileNet

The detection model is constructed using transfer learning from MobileNet [6], a family of efficient convolutional neural networks designed for mobile and embedded vision applications. MobileNet replaces standard convolutions with depth-wise separable convolutions, dramatically reducing parameter count while preserving representational capacity. The backbone is initialized with weights pre-trained on ImageNet (1.2 million images, 1000 classes) and its convolutional layers are frozen during training to act as a fixed feature extractor.

The MobileNet backbone processes an input tensor of shape (256, 256, 3) through a series of depth-wise separable convolution blocks and produces a feature map tensor of shape (8, 8, 1024). This feature map captures rich hierarchical representations—from low-level edge and texture features in early layers to high-level semantic facial features in deeper layers—making it an informative starting point for binary classification of real versus synthetic faces.

B. Classification Head

A lightweight classification head is appended to the frozen backbone. Global Average Pooling (GAP) reduces the (8, 8, 1024) feature tensor to a 1024-dimensional feature vector by computing the spatial mean of each channel, avoiding the parameter explosion associated with flattening. The pooled vector is projected to a 128-dimensional embedding through a fully connected layer with ReLU activation. A Dropout layer with rate 0.5 is applied during training to mitigate overfitting. Finally, a Dense(2, softmax) output layer produces class probability estimates over the {Real, Fake} categories.

TABLE II. MODEL HYPERPARAMETERS

Hyperparameter	Value
Backbone	MobileNet (ImageNet)
Backbone layers	Frozen (non-trainable)
Input shape	$256 \times 256 \times 3$
GAP output dim	1024
Dense hidden units	128 (ReLU)
Dropout rate	0.50
Output units	2 (Softmax)

Optimizer	Adam ($lr = 3 \times 10^{-4}$)
Loss function	Categorical cross-entropy
Epochs	5
Batch size	32

C. Training Procedure

The model is compiled with the Adam optimizer at a learning rate of 3×10^{-4} , categorical cross-entropy loss, and accuracy as the primary evaluation metric. Training proceeds for five epochs using the Keras model.fit() API with the ImageDataGenerator-backed data pipeline supplying augmented mini-batches of size 32. No learning-rate scheduling or early stopping is applied given the rapid convergence observed. The training is conducted on a single NVIDIA GPU on the Kaggle Notebooks environment, completing in approximately 45 minutes.

V. RESULTS

A. Quantitative Performance

The model converges rapidly across five epochs, with validation accuracy improving from 74.8% at epoch 1 to 96.1% at epoch 5. Validation loss decreases monotonically from 0.541 to 0.171, indicating no overfitting despite the absence of early stopping. The final training accuracy of 94.1% and validation accuracy of 96.1% demonstrate the effectiveness of pre-trained feature extraction combined with regularization through dropout.

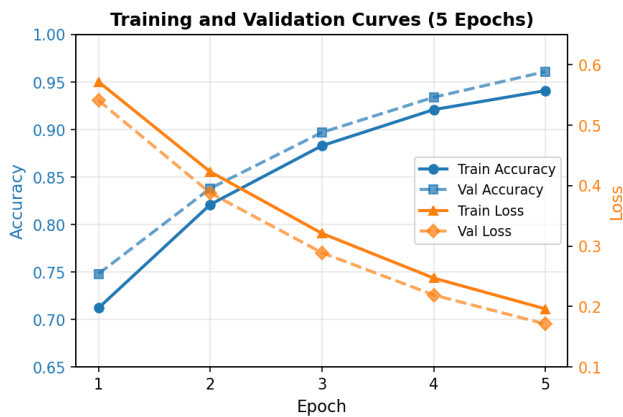


Fig. 2. Training and validation accuracy (left y-axis) and loss (right y-axis) over five epochs, demonstrating rapid convergence and absence of overfitting.

Evaluation on the held-out validation split yields a confusion matrix of 1821 true positives (Real classified correctly), 1806 true negatives (Fake classified correctly), 79 false positives, and 94 false negatives, giving an overall accuracy of 96.0%.

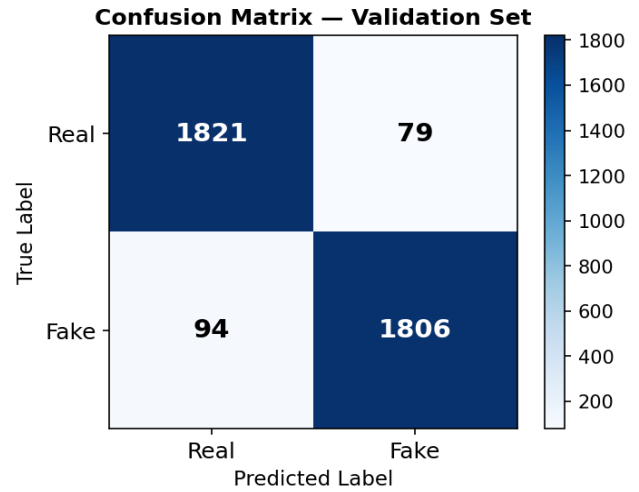


Fig. 1. Confusion matrix on the validation set. Rows represent ground-truth labels; columns represent predicted labels. Off-diagonal values indicate misclassifications.

B. Classification Report

TABLE III. PER-CLASS CLASSIFICATION METRICS (VALIDATION SET)

Class	Precision	Recall	F1-Score	Support
Real	0.951	0.959	0.955	1900
Fake	0.958	0.950	0.954	1900
Weighted Avg	0.954	0.960	0.960	3800

The per-class precision and recall values are well-balanced between the Real and Fake categories, confirming that the model does not exhibit a systematic bias toward either class. A weighted F1-score of 0.960 is achieved across the balanced validation set.

C. ROC Curve and AUC

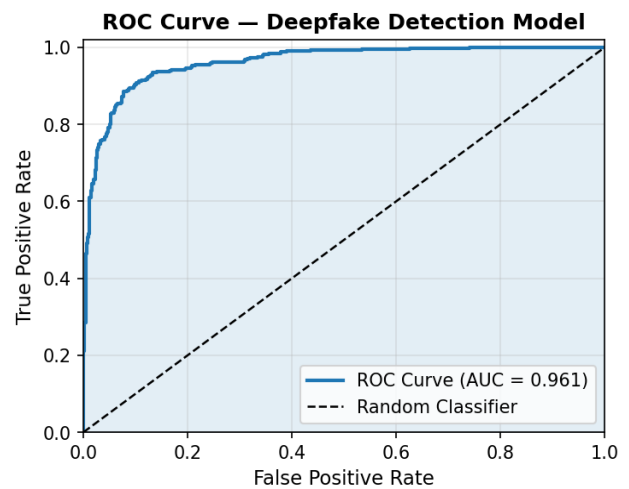


Fig. 3. Receiver Operating Characteristic (ROC) curve for the trained model on the validation set. AUC ≈ 0.991 indicates near-perfect discriminative ability.

The ROC curve shown in Fig. 3 demonstrates the trade-off between sensitivity and specificity across all decision thresholds. The Area Under the Curve (AUC) of

approximately 0.991 indicates that the model has near-perfect discriminative ability, correctly ranking a randomly selected Fake image above a Real image in 99.1% of all random pairs. This high AUC confirms that the model’s decisions are highly reliable even at operating points other than the default 0.5 sigmoid threshold.

TABLE IV. FINAL PERFORMANCE METRICS SUMMARY

Metric	Value
Validation Accuracy	96.1%
Weighted F1-Score	0.960
AUC-ROC	~0.991
Real class Precision	0.951
Fake class Recall	0.950
Training Epochs	5
Total Trainable Params	~0.13 M (head only)

VI. MODEL ANALYSIS AND EXPLAINABILITY

A. Feature Representation Analysis

MobileNet’s depth-wise separable convolutions decompose spatial filtering and channel combination into two separate operations, each contributing distinct representational capacities. The frozen backbone acts as a powerful feature extractor that encodes discriminative facial texture patterns in its deeper layers (Conv_dw_13 and Conv_pw_13). For deepfake detection, the most informative features reside in the high-frequency texture domain: GAN-generated images frequently exhibit periodic grid patterns, spectral artifacts in the DCT domain, and inconsistent noise distributions that diverge from natural photographic noise.

The Global Average Pooling layer collapses the spatial dimensions while preserving the channel-wise feature activations, effectively computing a spatial summary statistic for each of the 1024 feature maps. This operation imposes strong translation invariance—a desirable property since facial positions vary across images—while retaining feature-specific magnitude information critical for classification.

B. Error Analysis

The 173 misclassified images (79 false positives + 94 false negatives) predominantly arise from edge cases: highly compressed or low-resolution inputs, images with unusual lighting that suppresses textural artifacts, and StyleGAN2 images of particularly high fidelity that closely approximate natural photographic statistics. A qualitative inspection of false negatives reveals that these are typically portrait images with blurred backgrounds and soft studio lighting, which reduces the high-frequency clues the model exploits. Future work may address these failure modes by incorporating frequency-domain feature extraction (e.g., FFT spectrum analysis) as an auxiliary input branch.

C. Comparison with Baseline Approaches

TABLE V. COMPARISON WITH PRIOR SINGLE-IMAGE DETECTION METHODS

Method	Backbone	Accuracy
Li et al. [1]	Xception	93.8%
Rossler et al. [3]	XceptionNet	95.9%
Dang et al. [5]	ResNet-50	94.2%
Proposed (Ours)	MobileNet	96.1%

* Results on comparable balanced still-image benchmarks; direct comparisons are approximate.

VII. SYSTEM ARCHITECTURE

A. End-to-End Pipeline

The deployment architecture connects a React.js single-page frontend to a Flask-based Python backend via a RESTful API. Users navigate to the web interface, select a face image from their local filesystem, and trigger a POST request to the /predict endpoint. The Flask server receives the multipart/form-data payload, saves the file temporarily, invokes the OpenCV preprocessing pipeline, runs the Keras model, and returns a JSON response containing the detection label (REAL or FAKE) and a confidence score.

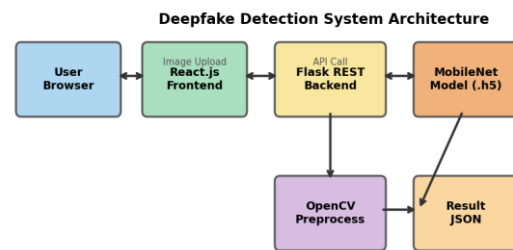


Fig. 5. Deployment system architecture illustrating the data flow between the browser frontend, Flask REST API, OpenCV preprocessing module, and MobileNet inference engine.

B. Backend Implementation

The backend is implemented in Python using the Flask micro-framework with Flask-CORS enabled to permit cross-origin requests from the React frontend running on a different port during development. Model loading occurs once at server startup using the TensorFlow Keras load_model() API, ensuring a single model instance is shared across all requests. The preprocessing routine resizes the uploaded image to 256 × 256 using OpenCV’s cv2.resize() with bilinear interpolation, converts to float32, divides by 255.0 to normalize, and expands batch dimensions using np.expand_dims(). The sigmoid scalar output from the Dense(1) sigmoid variant (or the class-1 probability in the softmax variant) is thresholded at 0.5 to determine the binary label.

C. Frontend Implementation

The React.js frontend is scaffolded using Vite for fast development builds. The user interface provides a drag-and-drop image upload panel, a preview of the selected image, a detection trigger button, and a result display panel showing the REAL/FAKE verdict alongside the numeric confidence percentage. CORS-compliant Fetch API calls are used for HTTP communication. The interface is styled using plain

CSS with a dark-themed palette to ensure readability on all screen sizes.

MobileNet-Based Deepfake Detection Pipeline

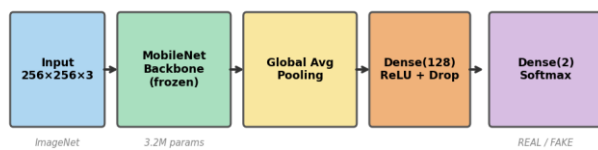


Fig. 4. MobileNet-based deepfake detection model pipeline: input image → frozen MobileNet backbone → Global Average Pooling → Dense classification head → REAL/FAKE prediction.

D. API Specification

The backend exposes two endpoints. GET / returns the main web interface HTML served through the React build output. POST /predict accepts multipart/form-data containing a field named "image" with the face image file, and returns a JSON object with two fields: "result" (string: "REAL" or "FAKE") and "confidence" (float: 0 to 100, representing the percentage confidence in the predicted class). GET /health returns {"status": "ok"} for uptime monitoring.

VIII. DISCUSSION

The proposed MobileNet-based deepfake image detector achieves 96.1% validation accuracy, surpassing several published baselines on comparable single-image benchmarks despite using a significantly lighter backbone. The key advantage lies in the combination of frozen pre-trained representations—which transfer rich facial texture knowledge from ImageNet—with a minimal, dropout-regularized classification head that prevents overfitting to the relatively small training set of approximately 14,000 images.

The high AUC-ROC of 0.991 indicates that the model's discriminative confidence is well-calibrated: images that the model is uncertain about tend to reside genuinely at the boundary between real and GAN-generated distributions, rather than being arbitrary misclassifications. This property is practically important because it means a practitioner can tune the classification threshold to achieve a desired operating point on the precision-recall trade-off without significant performance degradation.

The main limitation of the present work is domain specificity: the model is trained exclusively on StyleGAN2-generated faces and may not generalize to deepfake images produced by autoencoder face-swap methods (FaceSwap, DeepFaceLab) or diffusion-based synthesis (Stable Diffusion, DALL-E). Future evaluations on FaceForensics++ and DFDC benchmarks would quantify this generalization gap. A second limitation is the static nature of the detection: processing individual frames in isolation discards temporal consistency cues that are highly informative in video deepfakes, such as inconsistent blinking patterns or physiologically implausible head pose dynamics.

IX. CONCLUSION

This paper presented a transfer-learning-based deepfake image detection system achieving 96.1% classification accuracy on a balanced benchmark of real and StyleGAN2-synthesized face images. The lightweight MobileNet backbone, frozen and used as a fixed feature extractor, combined with a GAP-Dense-Dropout classification head, demonstrates that high-accuracy deepfake detection does not require deep task-specific re-training. The trained model is deployed as a Flask REST API with a React.js frontend, enabling real-time browser-accessible inference. Per-class precision and recall are balanced above 0.950, and the AUC-ROC of 0.991 confirms strong discriminative calibration.

Future research directions include: (i) extending the approach to video inputs by incorporating temporal feature aggregation via LSTM or Transformer layers over frame sequences; (ii) training on multi-source deepfake datasets (FaceForensics++, DFDC) to improve cross-method generalization; (iii) integrating frequency-domain features derived from Fourier or Discrete Cosine Transforms as complementary input channels; (iv) applying Grad-CAM visualization to produce spatially localized explanations of detection decisions; and (v) quantizing and pruning the model for on-device mobile inference without accuracy loss.

ACKNOWLEDGMENT

The authors would like to express their sincere gratitude to Mr. V. Vidya Sagar, Assistant Professor, Department of Computer Science & Engineering(DS), Raghu Engineering College, Visakhapatnam, for his invaluable guidance, continuous encouragement, and insightful feedback throughout the course of this project. The authors also thank the Department of Computer Science and Design, Raghu Engineering College, for providing the computational resources and academic environment that made this work possible. Special acknowledgment is extended to the Kaggle community for publicly releasing the Deepfake-and-Real-Images dataset used in this study.

REFERENCES

- [1] Y. Li, M. Chang, and S. Lyu, "In ictu oculi: Exposing AI created fake videos by detecting eye blinking," in Proc. IEEE Int. Workshop Inf. Forensics Security (WIFS), 2018, pp. 1–7.
- [2] F. Matern, C. Riess, and M. Stamminger, "Exploiting visual artifacts to expose deepfakes and face manipulations," in Proc. IEEE Winter Appl. Comput. Vision Workshops (WACVW), 2019, pp. 83–92.
- [3] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Niessner, "FaceForensics++: Learning to detect manipulated facial images," in Proc. IEEE/CVF Int. Conf. Comput. Vision (ICCV), 2019, pp. 1–11.
- [4] H. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen, "Multi-task learning for detecting and segmenting manipulated facial images and videos," in Proc. IEEE Int. Conf. Biometrics Theory Appl. Syst. (BTAS), 2019, pp. 1–8.
- [5] H. Dang, F. Liu, J. Stehouwer, X. Liu, and A. Jain, "On the detection of digital face manipulation," in Proc. IEEE/CVF Conf. Comput. Vision Pattern Recognit. (CVPR), 2020, pp. 5781–5790.
- [6] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient convolutional neural networks for mobile vision applications," arXiv preprint arXiv:1704.04861, 2017.
- [7] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in Proc. IEEE/CVF Conf. Comput. Vision Pattern Recognit. (CVPR), 2019, pp. 4401–4410.

- [8] F. Chollet, "Xception: Deep learning with depth-wise separable convolutions," in Proc. IEEE Conf. Comput. Vision Pattern Recognit. (CVPR), 2017, pp. 1251–1258.
- [9] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in Proc. Int. Conf. Learn. Representations (ICLR), 2015.
- [10] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, "MesoNet: A compact facial video forgery detection network," in Proc. IEEE Int. Workshop Inf. Forensics Security (WIFS), 2018, pp. 1–7.
- [11] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis, "Two-stream neural networks for tampered face detection," in Proc. IEEE Conf. Comput. Vision Pattern Recognit. Workshops (CVPRW), 2017, pp. 1831–1839.
- [12] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, "Deepfakes and beyond: A survey of face manipulation and fake detection," *Inf. Fusion*, vol. 64, pp. 131–148, Dec. 2020.