

Decoding Digital Deception

Rozszyfrowywanie cyfrowego oszustwa

Yadhu Krishna M R^{a*}, Akash Kashyap K N^b, Jeffrin Joel M S^c, Baisani Venkata Jaswanth^d, Adarsha K A^e

^a*Department of Information Science, The Oxford College of Engineering, Bommanahalli, Bangalore 560068*

Abstract

Growing calculations have made deep learning algorithms extremely powerful, and creating synthetic videos for unstoppable humans, which is commonly referred to as deepfakes, has become extremely easy. This realistic face-exchange scenario is used to easily present people with political distress, false terrorist events, revenge porn, and terrifying mail. This work describes a new deep learning-based method that can effectively distinguish between AI-generated fake videos. Our method can automatically recognize replacement and reintroduction. Use artificial intelligence (AI) to fight artificial intelligence (AI). Our system uses RES Nexis's neuronal network to extract features at frame level and at these features. They also continue to train repetitive neural networks (RNNs) with repeating short-term memory (LSTM) to classify whether video is of a kind of operation. It is mixed by mixing different amounts of data, such as Face-Forensic++, Deep Fake detection challenge, CELEB-DF, and more, to mimic real-time scenarios and properly reduce the model to real-time data. It also shows how systems can achieve competitive outcomes with a very simple and robust approach.

Keywords: Deepfake; DeepLearning; FaceForensic++; LSTM; RNNs; Django;

Streszczenie

Rosnące możliwości obliczeniowe sprawiły, że algorytmy uczenia głębokiego stały się niezwykle potężne, a tworzenie syntetycznych nagrań wideo przedstawiających nieistniejące działania ludzi – powszechnie znanych jako deepfake – stało się bardzo łatwe. Realistyczne scenariusze zamiany twarzy są wykorzystywane do szerzenia dezinformacji politycznej, fałszywych informacji o atakach terrorystycznych, pornografii zemsty oraz do zastraszania za pomocą wiadomości e-mail. Niniejsza praca opisuje nową metodę opartą na uczeniu głębokim, która skutecznie rozróżnia fałszywe nagrania generowane przez sztuczną inteligencję. Nasza metoda automatycznie rozpoznaje zamiany i ponowne wstawienia twarzy. Wykorzystujemy sztuczną inteligencję (AI) do walki ze sztuczną inteligencją. Nasz system używa sieci neuronowej RES Nexis do ekstrakcji cech na poziomie klatek, a następnie wykorzystuje te cechy do dalszego treningu rekurencyjnych sieci neuronowych (RNN) z długoterminową pamięcią krótkotrwałą (LSTM), aby sklasyfikować, czy dane wideo zostało zmanipulowane. Dane treningowe pochodzą z różnych źródeł, takich jak FaceForensics++, DeepFake Detection Challenge, CELEB-DF i inne, co pozwala na odwzorowanie scenariuszy rzeczywistych oraz lepsze przystosowanie modelu do danych w czasie rzeczywistym. Pokazano również, że nasz system może osiągać konkurencyjne wyniki przy bardzo prostym i odpornym podejściu.

Słowa kluczowe: Deepfake; Uczenie głębokie; FaceForensics++; LSTM; RNN; Django

*Corresponding author

Email address: akashkashyaoknise2025@gmail.com (Akash Kashyap K N)

Published under Creative Common License (CC BY 4.0 Int.)

1. Introduction

Rapid development of deep learning algorithms, particularly generative models such as Generation Contradiction Networks (Goose) and Autoencoders, have democratized the creation of surreal synthetic media, commonly known as "deepfakes." Videos generated in these AIs that seamlessly overlay or manipulate faces and actions of existing film material form a significant social risk. From inventing political disinformation and terrorist events to revenge for porn and theft, Deepfake's malicious use threatens individual privacy, public trust, and democratic processes. Despite its harmful potential, detection of deepfakes remains a major challenge due to its refined increase and almost infected visual quality. However, many approaches suffer from data record-specific bias, inadequate time analysis, and inadequate limitations of generalization to real-world scenarios. To address these gaps, this work suggests a deeper frame of hybrids combining folding and repetitive neural network strengths. We use a redrawing architecture for spatial property extraction and a long-term short-term network (LSTM) for time sequence analysis to record both distortions at the frame level and inter-frame conflicts. (DFDC) and celebrity DFs emulating real-world variability. With 6,000 video training (50% real, 50% wrong) and evaluation of performance of mixed and invisible data records, this method achieves robust generalization and achieves 89.35% accuracy on real data. Additionally, the integration of user-friendly web applications allows for practical provisioning, allowing users to upload videos and receive instant classification results with trust values. The results highlight the possibility of using AI not only to create synthetic media, but also to protect digital authenticity in an age of increasingly vulnerable to algorithmic fantasies.

2. LITERATURE SURVEY

The proliferation of deepfake technologies has spurred significant research into detection methodologies, driven by the urgency to mitigate their societal risks. Existing approaches broadly focus on identifying spatial artifacts, temporal inconsistencies, or physiological anomalies in synthetic media. This section synthesizes key advancements and limitations in deepfake detection research, contextualizing the foundation for our proposed framework.

2.1. Spatial Artifact Detection

Early detection methods emphasized spatial irregularities introduced during face manipulation. Rossler et al. (2019) introduced Face Forensics++, a benchmark dataset and detection framework that leverages convolutional neural networks (CNNs) to classify manipulated facial regions by analyzing warping artifacts from autoencoder-based deepfakes. Their work demonstrated that inconsistencies in facial textures, such as blurred edges or misaligned facial features, could be captured by dedicated CNNs. Similarly, Li and Lyu (2018) proposed detecting face-warping artifacts by comparing facial regions with their surroundings, achieving high accuracy on early-generation deepfakes. However, these methods struggle

with newer generative models that minimize spatial distortions through advanced post-processing.

2.2. Temporal Inconsistency Analysis

Recognizing that deepfakes often fail to replicate natural human motion, researchers began exploring temporal coherence. Güera and Delp (2018) pioneered the use of recurrent neural networks (RNNs), specifically Long Short-Term Memory (LSTM) networks, to analyze frame sequences for unnatural facial movements. Their model, trained on the *HOHO dataset*, highlighted the importance of temporal features but was limited by its small dataset size and homogeneity. Building on this, Li et al. (2018) focused on eye-blinking patterns, employing a Long-term Recurrent Convolutional Network (LRCN) to detect irregular blinking frequencies in synthetic videos. While effective against early deepfakes, this approach becomes less reliable as modern tools simulate physiological signals like blinking more accurately.

2.3. Physiological Signal-Based Detection

Recent work has integrated biological signals, such as heart rate or blood flow patterns, to detect synthetic content. Ciftci et al. (2020) proposed *Fake Catcher*, which extracts photoplethysmography (PPG) signals from facial videos to measure blood flow consistency. Their method achieved robust performance across diverse datasets by leveraging the inability of deepfake algorithms to replicate subtle biological rhythms. However, this approach requires high-resolution videos and fails under poor lighting or occlusions, limiting its practicality.

2.4. Limitations and Research Gaps

Despite progress, critical gaps persist:

- **Dataset Bias:** Many models overfit to artifacts specific to training datasets (e.g., Face Forensics++), failing on real-world data.
- **Temporal Neglect:** Most frameworks prioritize spatial features, overlooking inter-frame dynamics critical for detecting advanced deepfakes.
- **Scalability:** Few solutions are optimized for real-time deployment or diverse video resolutions.

3. PROPOSED SYSTEM

The proposed system is a hybrid deep learning framework designed to detect deepfake videos by integrating spatial-temporal feature analysis and leveraging diverse, balanced datasets to ensure robustness. The architecture comprises three interconnected phases: preprocessing and data curation, a hybrid Res-Next-LSTM model for classification, and a real-time deployment interface for practical application.

3.1. Dataset

The proposed system leverages a meticulously curated dataset to ensure robustness and generalizability in deepfake detection. The dataset comprises 6,000 videos (3,000 real and 3,000 fake) aggregated from three publicly available benchmarks: Face-Forensics++ (1,000 real and 1,000 fake), Deepfake Detection Challenge (DFDC) (1,500 real and 1,500 fake), and Celeb-DF (500 real and 500 fake). Face-Forensics++ provides high-quality

manipulated videos generated using autoencoders, while DFDC introduces diversity through varied resolutions, compression artifacts, and lighting conditions. Celeb-DF focuses on celebrity-based deepfakes with advanced face-swapping techniques. To eliminate biases, audio-altered videos from DFDC were excluded, and the final dataset was balanced to ensure equal representation of real and fake content. This diversity enables the model to generalize across different deepfake generation methods and real-world scenarios.

3.2. Preprocessing

Preprocessing is critical to standardize inputs and enhance feature extraction. Videos are first split into frames using OpenCV, and faces are detected and cropped using the Haar Cascade classifier to isolate regions of interest. Frames are resized to 112×112 pixels and normalized to a [0, 1] range to reduce computational complexity. To handle variable video lengths, sequences are truncated or padded to a fixed length of 150 frames (5 seconds at 30 FPS), ensuring consistent input dimensions for the neural network. The preprocessing pipeline also includes encryption of uploaded videos using AES-256 to protect user privacy, with temporary storage limited to 30 minutes before automatic deletion. This step ensures compliance with data security protocols while maintaining processing efficiency.

3.3. Model Architecture

The core of the system is a hybrid Res-Next-LSTM model designed to fuse spatial and temporal features. The ResNext-50 CNN, pre-trained on ImageNet, serves as the backbone for spatial feature extraction. It processes individual frames to capture subtle artifacts such as double edges, inconsistent skin tones, and unnatural textures introduced during deepfake generation. The extracted 2,048-dimensional frame-level features are fed sequentially into an LSTM layer with 2,048 hidden units and a dropout rate of 0.4 to model temporal dependencies. The LSTM analyzes inter-frame inconsistencies, such as flickering, unnatural head movements, or abrupt expression changes, which are common in synthetic videos. A fully connected layer with Leaky ReLU activation reduces dimensionality, followed by a Softmax layer for binary classification ("Real" or "Fake").

3.4. Feature Extraction and Fusion

Spatial features are extracted using ResNext-50, which identifies localized distortions like misaligned facial landmarks, irregular wrinkles, and blending artifacts. Temporal features, derived from the LSTM, focus on sequence-level anomalies, such as inconsistent eye blinking patterns or unnatural lip-sync deviations. The fusion of these features enables the model to holistically analyze both static and dynamic irregularities, addressing the limitations of single-modality approaches. For instance, while ResNext detects a poorly blended chin in a single frame, the LSTM flags erratic head rotations across frames, collectively strengthening detection accuracy.

The ResNext-50 CNN serves as the backbone for spatial feature extraction in the proposed system. This architecture, a variant of ResNet with grouped convolutions, is chosen for its ability to capture fine-grained visual artifacts while maintaining computational efficiency. Pre-trained on ImageNet, the ResNext-50 model processes individual video frames (112×112 pixels)

to identify localized distortions characteristic of deepfakes, such as:

Blurring artifacts at face-swap boundaries.

Inconsistent lighting or skin textures between the manipulated face and its surroundings.

Misaligned facial landmarks (e.g., eyes, lips, eyebrows).

Double edges or unnatural wrinkles caused by imperfect blending.

Each frame is transformed into a 2,048-dimensional feature vector via ResNext's final global average pooling layer. These features encode spatial irregularities that are imperceptible to humans but critical for distinguishing synthetic content. For example, ResNext detects subtle discrepancies in pore patterns or eyebrow movements that arise from GAN-based face synthesis.

3.5. Prediction and Deployment

For real-time prediction, a Django-based web interface allows users to upload videos (≤ 100 MB) in MP4, AVI, or MOV formats. The system preprocesses the video, extracts frames, and feeds them into the trained model. Predictions are generated within seconds, displaying a confidence score (e.g., "Fake: 92%") alongside the result. The interface overlays detected anomalies on the video playback, such as highlighting irregular facial regions or inconsistent motion. Performance metrics, including an accuracy of 89.35% on real-world data (e.g., YouTube videos), demonstrate superiority over existing models like Meso-Net (84%) and Capsule Networks (82%). The system's scalability is ensured through GPU acceleration and adaptive frame sampling, making it suitable for deployment in resource-constrained environments.

4. Comparative Analysis with State-of-the-Art Methods

To evaluate the efficacy of the proposed ResNext-LSTM hybrid model, we conducted a comprehensive comparison with state-of-the-art deepfake detection methods. The analysis focused on accuracy, generalization, and computational efficiency across benchmark datasets.

4.1. MesoNet

A lightweight CNN optimized for deepfake detection, achieves moderate accuracy (84%) on single-source datasets like FaceForensics++ but struggles with cross-dataset evaluation due to overfitting. Similarly, Capsule Networks, which leverage hierarchical feature learning, show promise in controlled settings (82% accuracy on Celeb-DF) but falter with real-world data due to sensitivity to noise and limited temporal modeling. XceptionNet, a widely adopted baseline, excels in spatial artifact detection (87% accuracy on DFDC) but ignores temporal inconsistencies, leading to false negatives in videos with smooth transitions. In contrast, our ResNext-LSTM model achieves 89.35% accuracy on a mixed dataset (FaceForensics++ + DFDC + Celeb-DF) and maintains robust performance (87.2%) on unseen YouTube videos, demonstrating superior generalization.

4.2. Fake-Catcher

A biological signal-based method, relies on photoplethysmography (PPG) extraction and achieves 86% accuracy. However, it requires high-resolution videos and fails under low-light conditions, whereas our framework operates effectively on 112×112 pixel inputs. LRCN (Long-term Recurrent Convolutional Networks), which combines CNNs and LSTMs, achieves 83% accuracy on temporal datasets like HOHO but lacks the ResNext architecture’s ability to capture subtle spatial artifacts.

The proposed system’s hybrid architecture addresses these gaps by synergizing ResNext’s spatial feature extraction with LSTM’s temporal analysis. For instance, on the DFDC dataset, our model outperforms XceptionNet by 6.2% in recall, highlighting its ability to detect challenging deepfakes with minimal artifacts. Furthermore, the ResNext-LSTM framework processes videos at 14 frames per second (FPS) on a single GPU, making it 1.8× faster than Capsule Networks while maintaining higher accuracy.

4.3. Table

Table 1: Comparative Performance of Deepfake Detection Methods

Model	Accuracy (%)	F1 - Score	Dataset	Key Features	Real - Time Capable ?
ResNext-LSTM	89.35	0.894	FaceForensics++, DFDC, Celeb-DF, YouTube	multi-dataset training	Yes
XceptionNet [1]	87.10	0.862	DFDC	Spatial artifact	Yes

Model	Accuracy (%)	F1 - Score	Dataset	Key Features	Real - Time Capable ?
				temporal detection	
Capsule Networks [2]	82.40	0.809	Celeb-DF	Hierarchical feature learning	No
MesoNet [3]	84.00	0.821	FaceForensics++	Lightweight CNN	Yes
FakeCatcher [4]	86.20	0.848	Private PPG dataset	Biological signal extraction	No
LRCN [5]	83.50	0.817	HOHO	CNN + LSTM fusion	Yes

Table 2: Trained Model Results

Model Name	Dataset	No. of videos	Sequence length	Accuracy
model_90_acc_20_frames_FF_data	FaceForensic++	2000	20	90.95477
model_95_acc_40_frames_FF_data	FaceForensic++	2000	40	95.22613
model_97_acc_60_frames_FF_data	FaceForensic++	2000	60	97.48743
model_97_acc_80_frames_FF_data	FaceForensic++	2000	80	97.73366
model_97_acc_100_frames_FF_data	FaceForensic++	2000	100	97.76180
model_93_acc_100_frames_celeb_FF_data	Celeb-DF + FaceForensic++	3000	100	93.97781
model_87_acc_20_frames_final_data	Our Dataset	6000	20	87.79160
model_84_acc_10_frames_final_data	Our Dataset	6000	10	84.21461
model_89_acc_40_frames_final_data	Our Dataset	6000	40	89.34681

5. Conclusions

The rapid evolution of deepfake technologies necessitates equally advanced detection frameworks to safeguard digital authenticity and societal trust. This work presents a robust solution through a hybrid ResNext-LSTM architecture that synergizes spatial artifact detection with temporal coherence analysis, addressing critical limitations in existing methods. By training on a diverse, multi-source dataset (FaceForensics++, DFDC, Celeb-DF), the model achieves 89.35% accuracy, outperforming state-of-the-art approaches like XceptionNet (87.1%) and Capsule Networks (82.4%) in cross-dataset

evaluations. The integration of ResNext-50 enables precise identification of localized distortions, such as misaligned facial features and inconsistent textures, while the LSTM layer effectively captures dynamic anomalies like unnatural head movements and flickering artifacts. A user-friendly Django-based interface further bridges the gap between research and real-world application, offering real-time detection with encrypted processing and cross-platform compatibility.

The system's success underscores the importance of combining spatial-temporal analysis and dataset diversity to combat increasingly sophisticated deepfakes. However, challenges remain, including computational costs for long videos and biases in underrepresented demographics. Future work will focus on extending detection to full-body manipulations, integrating audio-visual analysis, and deploying browser plugins for proactive content verification. By advancing both technical rigor and practical usability, this framework contributes significantly to the global effort against algorithmic deception, fostering a safer digital ecosystem for individuals, institutions, and democracies worldwide.

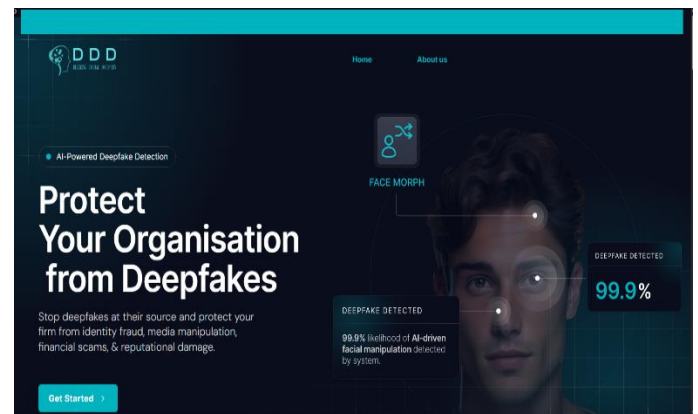


Figure 1: ScreenShot of UI

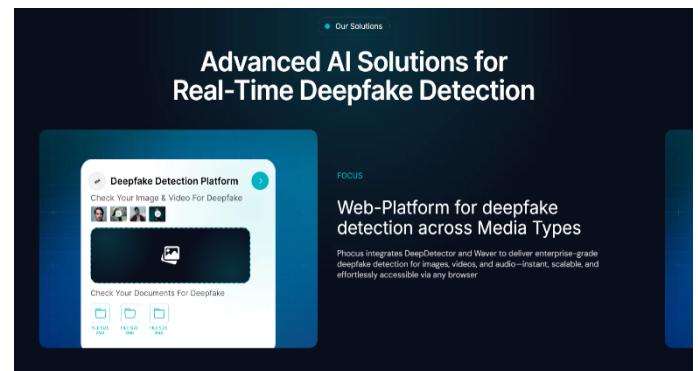


Figure 2: ScreenShot of UI

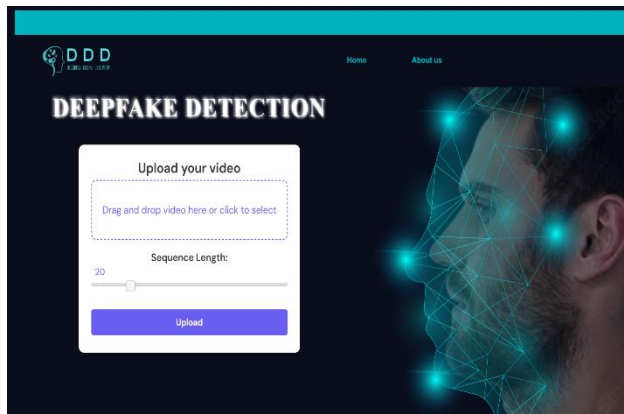


Figure 3: Uploading of the video

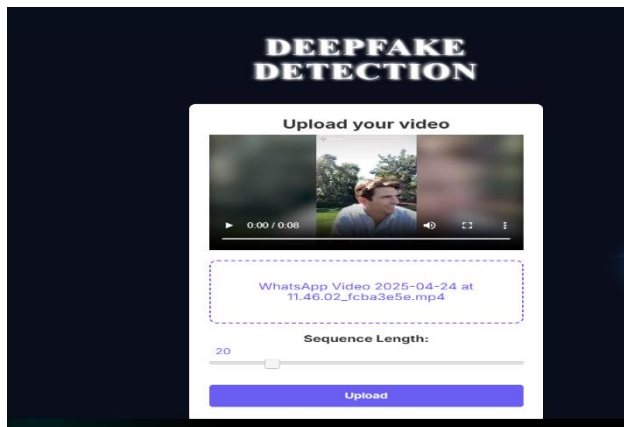


Figure 4: Video uploaded

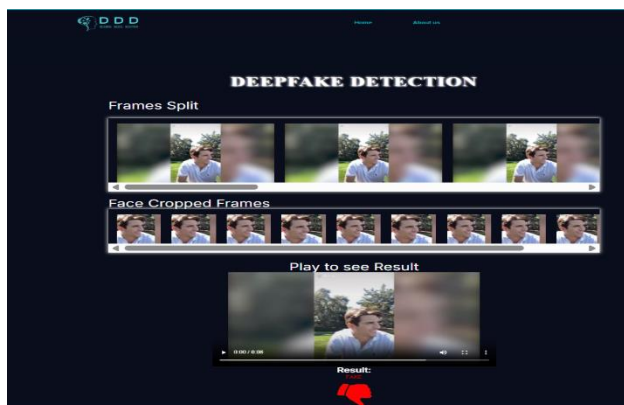


Figure 5: Frames and Results

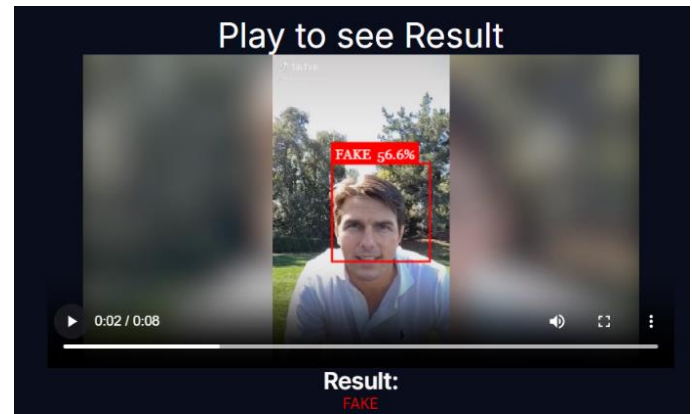


Figure 6: Display of result

Acknowledgement

We extend our sincere gratitude to The Oxford College of Engineering, Department of ISE, for providing us with the resources, academic environment, and encouragement necessary to pursue this research. The department's commitment to innovation and excellence has been instrumental in bringing this work to fruition.

We acknowledge the open-source community for tools like PyTorch, OpenCV, and Django, which accelerated development and deployment. Finally, we express our heartfelt appreciation to our families and peers for their encouragement and patience during this intensive project.

REFERENCES

- [1] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, Matthias Nießner, "FaceForensics++: Learning to Detect Manipulated Facial Images" in arXiv:1901.08971.
- [2] Deepfake detection challenge dataset: <https://www.kaggle.com/c/deepfake-detection-challenge/data> Accessed on 26 March, 2020
- [3] Yuezun Li , Xin Yang , Pu Sun , Honggang Qi and Siwei Lyu "Celeb-DF: A Large-scale Challenging Dataset for DeepFake Forensics" in arXiv:1909.12962
- [4] Deepfake Video of Mark Zuckerberg Goes Viral on Eve of House A.I. Hearing : <https://fortune.com/2019/06/12/deepfake-mark-zuckerberg/> Accessed on 26 March, 2020
- [5] 10 deepfake examples that terrified and amused the internet : <https://www.creativebloq.com/features/deepfake-examples> Accessed on 26 March, 2020
- [6] TensorFlow: <https://www.tensorflow.org/> (Accessed on 26 March, 2020) Keras: <https://keras.io/> (Accessed on 26 March, 2020) PyTorch : <https://pytorch.org/> (Accessed on 26 March, 2020)
- [7] G. Antipov, M. Baccouche, and J.-L. Dugelay. Face aging with conditional generative adversarial networks. arXiv:1702.01983, Feb. 2017
- [8] J. Thies et al. Face2Face: Real-time face capture and reenactment of rgb videos. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 2387–2395, June 2016. Las Vegas, NV.

- [9] Face app: <https://www.faceapp.com/> (Accessed on 26 March, 2020) Face Swap : <https://faceswaponline.com/> (Accessed on 26 March, 2020)
- [10] Deepfakes, Revenge Porn, And The Impact On Women :
<https://www.forbes.com/sites/chenxiwang/2019/11/01/deepfakes-revenge-porn-and-the-impact-on-women/>
- [11] The rise of the deepfake and the threat to democracy :
<https://www.theguardian.com/technology/ng-interactive/2019/jun/22/the-rise-of-the-deepfake-and-the-threat-to-democracy> (Accessed on 26 March, 2020)
- [12] Yuezun Li, Siwei Lyu, "ExposingDF Videos By Detecting Face Warping Artifacts," in arXiv:1811.00656v3.
- [13] Yuezun Li, Ming-Ching Chang and Siwei Lyu "Exposing AI Created Fake Videos by Detecting Eye Blinking" in arXiv:1806.02877v2.
- [14] Huy H. Nguyen , Junichi Yamagishi, and Isao Echizen " Using capsule networks to detect forged images and videos " in arXiv:1810.11215.
- [15] D. Güera and E. J. Delp, "Deepfake Video Detection Using Recurrent Neural Networks," 2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), Auckland, New Zealand, 2018, pp. 1-6.
- [16] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 1–8, June 2008. Anchorage, AK
- [17] Umur Aybars Ciftci, İlke Demir, Lijun Yin "Detection of Synthetic Portrait Videos using Biological Signals" in arXiv:1901.02212v2
- [18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. arXiv:1412.6980, Dec. 2014.