# Detection of Deep-Morphed Deepfake Images to Make Robust Automatic Facial Recognition Systems

Alakananda Mitra
Dept. of Computer Science and Engineering
University of North Texas, USA.
Email: AlakanandaMitra@my.unt.edu

Saraju P. Mohanty
Dept. of Computer Science and Engineering
University of North Texas, USA.
Email: saraju.mohanty@unt.edu

Peter Corcoran
School of Engineering and Informatics
National University of Ireland, Galway, Ireland.
Email: peter.corcoran@nuigalway.ie

Elias Kougianos
Dept. of Electrical Engineering
University of North Texas, USA.
Email: elias.kougianos@unt.edu

*Abstract*—Face Morphing has emerged as a pervasive attack of Facial Recognition Systems. The rapid growth of Generative Adversarial Networks takes it to a complete new level. Deepfake or deep neural network based face morphing, a.k.a deep-morph attack, presents a significant threat to Facial Recognition System. In this paper, we propose a novel Convolutional Neural Network based detection method of deep morphed deepfake images which is suitable for IoT environments in smart cities. A high accuracy of 94.83% has been achieved for the DeepfakeTIMIT HQ dataset. This lightweight and fast network is a natural choice for IoT environments.

*Index Terms*—Smart City, Facial recognition System, Deepfake, Deep-Morph, Deep Learning, Convolutional Neural Network.

## I. INTRODUCTION

The concept of smart urban society is changing the globe. In today's world, face recognition has given rise to biometric-based identification systems in use from law enforcement to traffic management, device security, crime prevention, or receiving access to any facility in smart cities [1] using advanced sensors, cameras and various IoT devices. Its wide use is mainly because it can be performed in a non-invasive way without touching any device. In the post COVID-19 world, this is a prerequisite condition.

Biometric based data is individual-specific, unique to the respective user, and easy to capture. Facial features are well-suited for any biometric based identification system as taking pictures is easy at any setting. Fig. 1 illustrates an application field of a facial recognition system (FRS) in a smart city. All the facilities of the smart city are accessed by citizens after verifying their identity through FRS.

But FRS is susceptible and vulnerable to various attacks such as presentation attacks, indirect/channel attacks, traditional face morphing attacks, deepfakes etc. Face Morphing Attacks (FMA), both traditional and deep morph, have arisen as serious contenders among these attacks. They fool the FRS by employing the same morphed photo for two people. The morphed image is generated by combining images of two
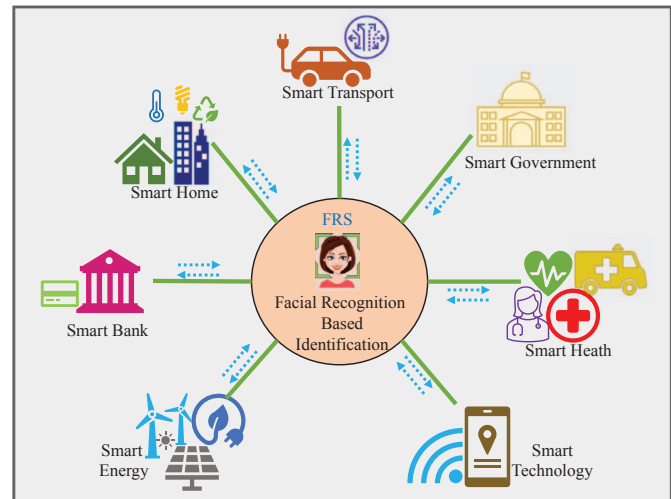


Fig. 1. An Example of Facial Recognition System (FRS) Application - Identification of Individuals in Smart City.

different people. As a result, the morphed image has the characteristics of both people to a certain degree. A certain level of similarity between two people makes this possible.

With the advent of deep neural networks, especially Generative Adversarial Networks (GANs) [2], face morphing has become easy and sophisticated. High quality deepfake images/videos are rampant in social media and in various websites. They change the perception of truth. These deep morph images pose threat to biometrics based facial recognition systems.

In this paper, we present a deep learning based detection method of deep morphed images in the context of a facial recognition system of smart cities. Here we use a Convolutional Neural Network (CNN), suitable for IoT devices, as feature extractor of the detection method.

The rest of the paper is organized into seven sections. Section II presents the issues of deep morphed images in the

context of FRS in smart cities along with an outline of the solution proposed. State-of-the-Art solutions are discussed in Section III. Our proposed work is presented in Section IV whereas its implementation is discussed in Section V. Results are stated in Section VI. Section VII concludes the paper with suggestions for future work.

## II. FACE MORPHING ATTACKS ON FRS OF SMART CITY

### A. Problem Addressed in the Current Paper

In this section, we discuss how face morphing attacks (FMA) affect facial recognition systems (FRS) of smart cities.

Fig. 2 illustrates face morphing attack. To make the facial recognition process of smart cities seamless and mass scale, the registration process should be easy.



(a) Enrollment with Existing Photo ID Containing Morphed Image
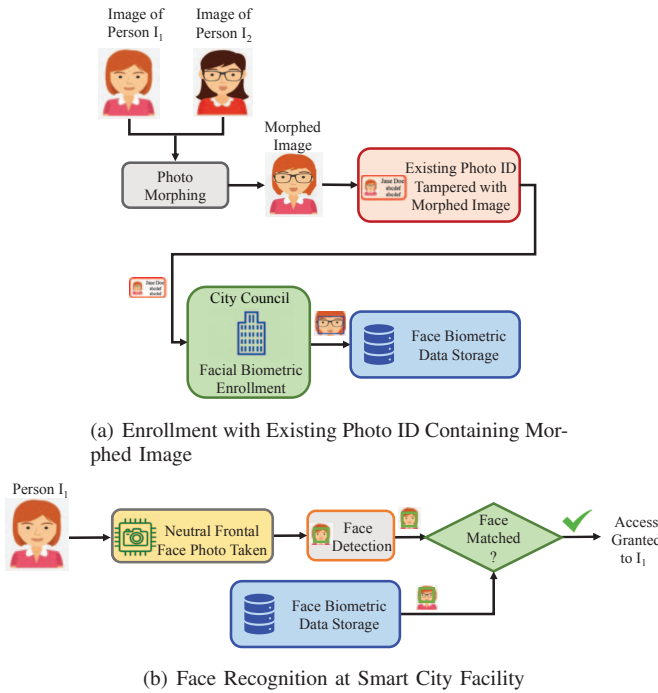


(b) Face Recognition at Smart City Facility

Fig. 2. Face Enrollment with Tampered ID and Recognition on FRS at Smart City

The suggested pipeline for registering a citizen in the FRS of that smart city is shown in Fig.3. During enrollment in the smart city FRS, citizens submit their existing photo id to the city council office. If the photo is matched with the person, he/she is registered in the city's FRS as in Fig. 2(a). But if this photo is morphed and the person is matched with the photo id as in Fig. 2(b), the facial recognition system will be used in the wrong way.



Fig. 3. Registration Pipeline

We assume that the submitted photo is morphed from hostile $I_1$'s and victim $I_2$'s photos and is used in the id of the hostile person $I_1$. Then the hostile person will be registered in the FRS with a photo id containing a *Morphed Image*. As a result, the hostile person $I_1$ will get access to smart city facilities by impersonating as in Fig. 2(b).

### B. Morphed Face Generation

An FMA can generate highly sophisticated morphed images and fool FRS. FMA is mainly of two types, depending on how they are generated.

*1) Landmark Based Morphing:* This is traditional face morphing. Landmark points of both faces are found and they are aligned by warping. Various techniques are followed to warp images. Finally, a blending operation is performed between the textures of both images. The similarity score of the morphed image depends on the blending parameter.

*2) DNN Based Morphing aka Deep Morph Deepfake:* Recent advances in GANs made the process of face morphing automated. GAN generated images look more sophisticated. In 2018, MorGAN [3] was introduced in the context of face morphing. It is a 3-module network consisting of an encoder, a decoder and a discriminator. The encoder and decoder form the generator which plays a zero sum game with the discriminator. Once the generator is trained, there is no need of the discriminator to generate the deep morph image of dimension $64 \times 64$. Images generated by StyleGAN [4], introduced in 2019, are of high resolution $1024 \times 1024$. They are visibly much better than MorGAN generated images. The recently introduced StyleGAN 2 [5] generates even higher quality morphed images than the original StyleGAN. It can transfer the expression of the source images. Deep morph deepfake images can also be generated with FSGAN [6]. No subject specific training is required for FSGAN to generate a deep morphed image. Face swapping results from FSGAN are much superior than existing deepfake face swap tools.

### C. The Solution Proposed

To address the problem in Section II-A, we propose a CNN based solution which can detect the deep morphed deepfake images submitted for registration in FRS of smart cities. We employ a CNN, well suited in an IoT enviroment, as the feature extractor and a fully connected layer as classifier. We achieve a high accuracy compared to existing works. We also study the features extracted by the CNN.

## III. RELATED WORKS

Facial forgery can be grouped into four categories: full face synthesis (e.g. StyleGAN [4], MorGAN [3]), identity swap (e.g. FSGAN [6]), attribute manipulation (e.g. StarGAN [7]) and expression transfer (e.g. Face2Face [8] [9]). As our point of interest is deepfake, here we discuss the detection techniques related to deepfake or face swapping.

We start our literature survey with an audio-visual feature based deepfake video detection method [10]. Principal Component Analysis (PCA) has been used along with different classification techniques. The DeepfakeTIMIT [11], [12] database

was used. Only visual aspects have been used in [13]. As a classifier, logistic regression and multi layer perceptron models have been utilized. The accuracy of the method largely varies with various datasets. [14] and [15] are also visual artifact based and are focused on deepfake detection. Higher accuracy with less computation makes those methods competitive. The proposed algorithms along with XceptionNet and classifier network make the method less compute intensive in detecting deepfake videos. An IoT friendly GAN generated deepfake image detection approach has been presented in [16]. Textural features of the images have been utilized to detect the fake images with a considerable accuracy.

Head pose along with facial expression has been used in [17]. dlib based facial landmark points detector and support vector machine have been used respectively to measure the head pose and to classify. Their model works well with UADFV dataset only. The same type of approach has been followed in [18] but with much higher accuracy. FSGAN [6] has been used to generate deepfake datasets.

Faceswap and deepfake based datasets have also been used in many papers. Face warping features are the base of detection in [19] whereas mesoscopic features are used in [20]. Overall, higher accuracy is obtained in [19] than [20]. A fusion model of steganalysis and deep learning features has been applied in [21] but accuracy of this model is lower than the previous two.

A capsule network based detection method has been proposed in [22]. The model has only higher accuracy in Face Forensics ++ dataset [23]. It shows lower accuracy in other tested datasets.

The papers mentioned above have employed the Deepfake-TIMIT (DF-TIMIT) [12] dataset. Other than [19], the accuracy of each model is not high for DF-TIMIT. That was our motivation to work on this dataset after working on various deepfake datasets in our previous papers [14], [15].

## IV. CNN BASED DETECTION OF DEEP MORPHED DEEPFAKE ATTACKS IN SMART CITIES CONTEXT

### A. Data Processing:

Before training our network, we process the data as per our requirements. As full face images have been used in training, we detect and align the face first. Then they are cropped, resized to $224 \times 224$, and normalized.

### B. Detection Network

To detect the deep morphed deepfake images, a CNN based method as shown in Fig. 4 is used. An existing CNN, suitable for facial recognition system in IoT devices of smart cities, has been used as the feature extractor of the network. The used CNN is based on depthwise separable convolution.

Depthwise separable convolution leverages the depthwise and pointwise filters, by performing depthwise convolution prior to pointwise convolution. The cost of depthwise separable convolution is much smaller than standard convolution. Each depthwise convolution filter is applied on each input
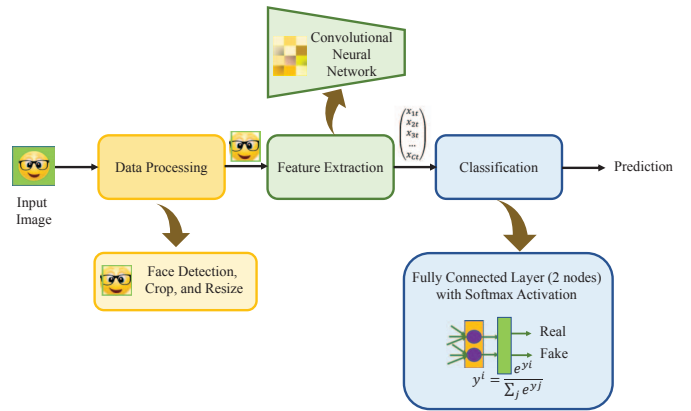


Fig. 4. Detection Method

channel or depth. A linear combination of depthwise filter output is computed by applying pointwise convolution filters.

MobileNetV2 [24] has been used as the feature extractor. It employs depthwise separable convolution. Here linear bottlenecks have been used between the layers. Shortcuts connect the bottleneck layers instead of non-bottleneck layers as in ResNets [25] and make it faster and accurate.

Tha last layer of the CNN is removed. A fully connected layer with *softmax* activation and 2 nodes has been used as the classifier.

## V. EXPERIMENTAL VALIDATION

### A. Dataset

Two public datasets have been used to evaluate the detection method. For fake images, we used the DeepfakeTIMIT (DF-TIMIT) [11], [12] dataset. The deepfake videos in this dataset have been generated using videos of 32 subjects from the VidTIMIT [26] dataset with FSGAN [6]. A total of 620 videos are generated using a lower quality (LQ) with $64 \times 64$ input/output sized model and a higher quality (HQ) $128 \times 128$ input/output sized model. We extracted fake image frames from those videos at a rate of $25 fps$.

For real images, the VidTIMIT [26] dataset has been used. It is a dataset with 43 individuals - 24 males and 19 females. The videos were shot in 3 different settings. There are 10 videos for each subject. The videos are stored as frames in the dataset. For our work, we used the same subject videos as DF-TIMIT [12]. Dataset details are mentioned in Table I.

TABLE I
DATASET DETAILS OF THE EXPERIMENT

| Real | | Fake | |
|------|------|------|------|
| Dataset Source | No. of Images | Dataset Source | No. of Images |
| VidTIMIT | 34,004 | DeepfakeTIMIT (HQ) | 33,988 |
| VidTIMIT | 34,004 | DeepfakeTIMIT (LQ) | 34,025 |

### B. Training Protocol

To train the network we follow the below protocol.

- Transfer learning has been used to reduce the training time. We chose an ImageNet trained MobileNetV2 [24] network as the feature extractor by replacing the last 1000 node fully connected layer with a 2 node fully connected layer. To set up the classifier layer, it is initially trained for 10 epochs keeping the weights of the feature extractor frozen. Then, the whole network has been trained for 15 more epochs from end-to-end. Finally the best model is chosen depending on the validation accuracy.
- Training data has been augmented to improve performance.
- The model has been trained on both datasets DF-TIMIT HQ and DF-TIMIT LQ separately. In both cases we use the same VidTIMIT dataset [26] for real images. During evaluation, the same dataset and cross dataset evaluation have been performed. For the same dataset evaluation, unseen data from the dataset is utilized for testing. For cross evaluation, one dataset is used for training and the other for testing. The results are presented in Sec. VI.
- For both datasets, ∼48,000 images for training, ∼12,000 for validation, and ∼4,000 images for testing have been used as shown in Table II.
- The network has been optimized with the Adam [27] optimizer of learning rate 0.0002 and other parameters to default values.

TABLE II
DATASET DIVISION FOR BOTH DF-TIMIT HQ AND LQ DATASET

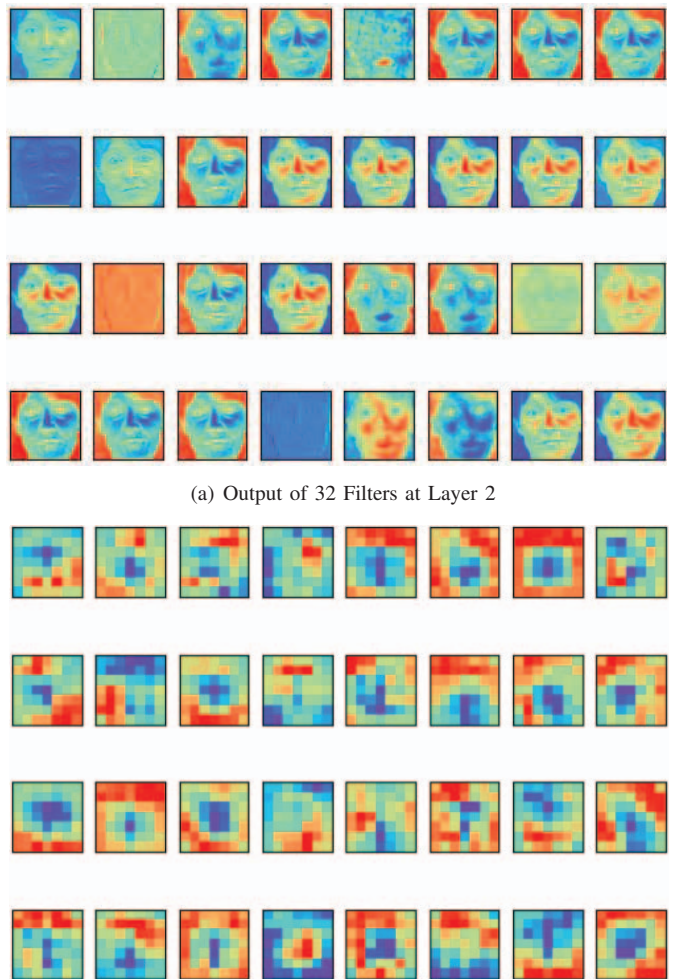| Data | Number of Images | | |
|---|---|---|---|
| | Real | Fake | |
| | | DF-TIMIT HQ | DF-TIMIT LQ |
| Train | 23,873 | 23,939 | 23965 |
| Validation | 6135 | 6000 | 6010 |
| Test | 3996 | 4049 | 4050 |

We implemented the detection method in Keras [28]. The model has been trained in a GeForce RTX 2060 laptop with 6GB shared and 16GB total memory.

## VI. RESULTS

Fig. 5 shows the features extracted by two different layers of MobileNetV2 [24]. Features extracted by layer 2 are shown in Fig. 5(a). All the filters of layer 2 are activated at each part of input. As we go deeper in the CNN, the layers extract more complex features. Features extracted by 32 filters out of 1280 filters of layer 153 are shown in Fig. 5(b).

Fig. 6(a) shows the class activation heatmap of feature extractor MobileNetV2 [24] pretrained on ImageNet using GRAD-CAM [29] and Fig. 6(b) is that of MobileNetV2 [24] trained on DF-TIMIT HQ dataset. Fig. 6 shows MobileNetV2 [24], trained on DF-TIMIT HQ, correctly classifies real and fake images whereas Imagenet [30] trained MobileNetV2 [24] fails to classify.
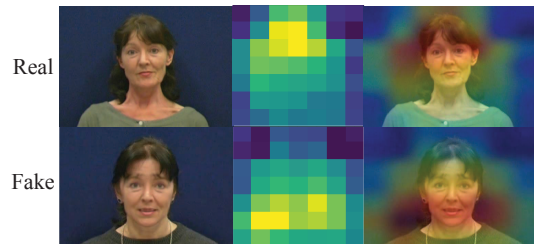
Table. III shows accuracy and inference times of different evaluation scenarios. The maximum accuracy has been obtained when the model is trained on DF-TIMIT LQ dataset and tested on the same dataset.
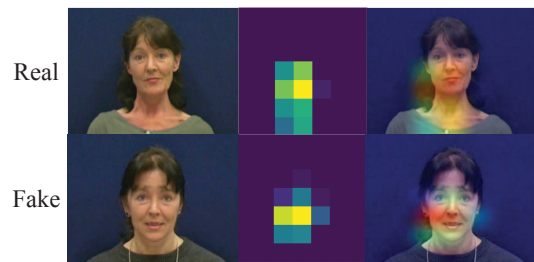


(a) Output of 32 Filters at Layer 2



(b) Output of 32 Filters Out of 1280 Filters at Layer 153

Fig. 5. Feature Visualization of MobileNetV2 for Sample Layers



(a) Pre-trained on ImageNet (Predicted Wrong)



(b) Trained on DF-TIMIT HQ Dataset (Predicted Correct)

Fig. 6. Class Activation Map Visualization using GRAD-CAM for MobileNetV2

| Training Dataset | Testing Dataset | Accuracy (%) | Inference Time (mS) |
|---|---|---|---|
| DF-TIMIT(HQ) | DF-TIMIT(HQ) | 94.83 | 3.67 |
| DF-TIMIT(LQ) | DF-TIMIT(LQ) | 100.00 | 3.76 |
| DF-TIMIT(HQ) | DF-TIMIT(LQ) | 96.91 | 3.81 |
| DF-TIMIT(LQ) | DF-TIMIT(HQ) | 57.38 | 4.45 |

* For real images → VidTIMIT dataset.

When the model is trained on the DF-TIMIT HQ dataset, the accuracy is high in both cases of testing. The accuracy is really poor when trained on low quality images and tested on high quality images, which is expected.

Inference time is also similar in the first three cases and high at the last case, as expected. When the model is trained on low quality images, it considers high quality fake images as real.

To evaluate classification performance of the model, the confusion matrix has been generated for a scenario when training and testing are performed on DF-TIMIT HQ dataset [12] as in Fig. 7. *Precision*, *Recall*, *Accuracy*, and *F1-score* have been calculated in Table IV from the confusion matrix.
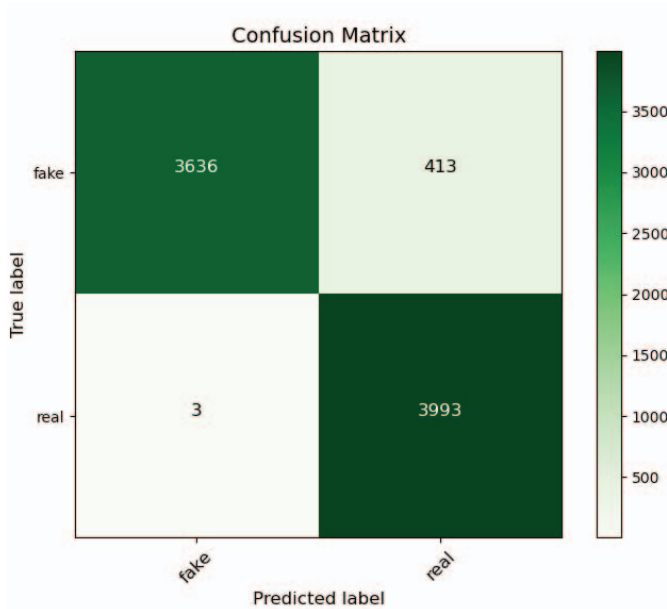


Fig. 7. Confusion Matrix - Trained and Tested on DF-TIMIT HQ

TABLE IV
CLASSIFICATION REPORT-TRAINED AND TESTED ON DF-TIMIT HQ

| Test images | Precision (%) | Recall (%) | F1-score |
|---|---|---|---|
| 3996 Real | 100.0 | 90.0 | 95.0 |
| 4048 Fake | 91.0 | 100.0 | 95.0 |
| Macro Average | 95.0 | 95.0 | 95.0 |
| Weighted Average | 95.0 | 95.0 | 95.0 |
| Total 8044 | Accuracy (%) | 95.0 | |

Table V compares our proposed method with other existing solutions. We achieve much higher accuracy using MobileNetV2 [24] as feature extractor and training the network with full face images.

TABLE V
PERFORMANCE COMPARISON OF OUR METHOD WITH
STATE-OF-THE-ART SOLUTIONS

| Study | Year | Performance(%) | |
|---|---|---|---|
| | | DF-TIMIT LQ | DF-TIMIT HQ |
| Matern et al. [13] | 2019 | AUC= 77.00 | AUC=77.30 |
| Yang et al. [17] | 2019 | AUC= 55.10 | AUC= 53.20 |
| Afchar et al. [20] | 2018 | AUC= 87.80 | AUC=68.40 |
| Zhou et al. [21] | 2018 | AUC= 83.50 | AUC=73.50 |
| Nguyen et al. [22] | 2019 | AUC= 78.40 | AUC=74.40 |
| **Proposed Method** | 2021 | ACC = 100.00 | ACC=94.83 |

*ACC → Accuracy ; *AUC → Area Under the Curve

## VII. CONCLUSION AND FUTURE WORK

As deep morphed deepfake images pose serious threat to biometric based facial recognition systems of smart cities, extensive research on deepfake video/image detection is needed to combat this problem. In this paper, a CNN based method has been proposed to detect FSGAN [6] generated deep morphed deepfake images. The lightweight model shows a lot of promise by achieving very high accuracy. As a future work, the model can be trained with more datasets for generalizability of the model. A one-shot learning approach might help in training the model.

## REFERENCES

[1] S. P. Mohanty, U. Choppali, and E. Kougianos, "Everything you wanted to know about smart cities: The internet of things is the backbone," *IEEE Consumer Electronics Magazine*, vol. 5, no. 3, pp. 60–70, 2016.

[2] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, Eds., vol. 27. Curran Associates, Inc., 2014, pp. 2672–2680.

[3] N. Damer, A. M. Saladié, A. Braun, and A. Kuijper, "Morgan: Recognition vulnerability and attack detectability of face morphing attacks created by generative adversarial network," in *Proceedings of IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, 2018, pp. 1–10.

[4] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4401–4410.

[5] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of stylegan," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020, pp. 8107–8116.

[6] Y. Nirkin, Y. Keller, and T. Hassner, "FSGAN: Subject agnostic face swapping and reenactment," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7184–7193.

[7] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo, "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8789–8797.

[8] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner, "Face2Face: Real-time Face Capture and Reenactment of RGB Videos," in *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2387–2395.

[9] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, "Deepfakes and beyond: A survey of face manipulation and fake detection," *Information Fusion*, vol. 64, pp. 131–148, 2020.

[10] Pavel Korshunov and S. Marcel, "Vulnerability of face recognition to deep morphing," *arXiv:abs/1910.01933*, p. 5 pages, 2019.

[11] P. Korshunov and S. Marcel, "Deepfakes: a new threat to face recognition? assessment and detection," *arXiv:abs/1812.08685*, 2018. [Online]. Available: http://arxiv.org/abs/1812.08685

[12] P. Korshunov and S. Marcel, "Deepfakes: a new threat to face recognition? assessment and detection," Idiap, Idiap-RR Idiap-RR-18-2018, 12 2018.

[13] F. Matern, C. Riess, and M. Stamminger, "Exploiting visual artifacts to expose deepfakes and face manipulations," in *Proceedings of IEEE Winter Applications of Computer Vision Workshops (WACVW)*, 2019, pp. 83–92.

[14] A. Mitra, S. P. Mohanty, P. Corcoran, and E. Kougianos, "A novel machine learning based method for deepfake video detection in social media," in *Proceedings of IEEE International Symposium on Smart Electronic Systems (iSES) (Formerly iNiS)*, 2020, pp. 91–96.

[15] A. Mitra, S. P. Mohanty, P. Corcoran, and E. Kougianos, "A machine learning based approach for deepfake detection in social media through key video frame extraction," *SN Computer Science*, vol. 2, no. 2, p. 98, 2021.

[16] A. Mitra, S. P. Mohanty, P. Corcoran, and E. Kougianos, "Easydeep: An iot friendly robust detection method for gan generated deepfake images in socialmedia," in *Proceedings of the 4th IFIP International Internet of Things (IoT) Conference (IFIP-IoT)*, 2021, Accepted, In Press.

[17] X. Yang, Y. Li, and S. Lyu, "Exposing deep fakes using inconsistent head poses," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 8261–8265.

[18] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, and H. Li, "Protecting world leaders against deep fakes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019, pp. 38–45.

[19] Y. Li and S. Lyu, "Exposing deepfake videos by detecting face warping artifacts," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019, pp. 46–52.

[20] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, "Mesonet: a compact facial video forgery detection network," in *Proceedings of IEEE International Workshop on Information Forensics and Security (WIFS)*, 2018, pp. 1–7.

[21] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis, "Two-stream neural networks for tampered face detection," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017, pp. 1831–1839.

[22] H. H. Nguyen, J. Yamagishi, and I. Echizen, "Use of a capsule network to detect fake images and videos," *arXiv:abs/1910.12467*, 2019. [Online]. Available: http://arxiv.org/abs/1910.12467

[23] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "FaceForensics++: Learning to detect manipulated facial images," in *Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 1–11.

[24] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

[25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[26] C. Sanderson and B. Lovell, "Multi-region probabilistic histograms for robust and scalable identity inference," Lecture Notes in Computer Science (LNCS), Vol. 5558, pp. 199–208, 2009.

[27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: http://arxiv.org/abs/1412.6980

[28] F. Chollet, "Keras," https://keras.io, 2015, Accessed on 15 July 2021.

[29] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618–626.

[30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of IEEE conference on computer vision and pattern recognition*, 2009, pp. 248–255.