



Deepfakes and beyond: A Survey of face manipulation and fake detection

Ruben Tolosana*, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, Javier Ortega-Garcia

Biometrics and Data Pattern Analytics Lab, Universidad Autonoma de Madrid, Spain

ARTICLE INFO

Keywords:

Fake news
Deepfakes
Media forensics
Face manipulation
Face recognition
Benchmark
Databases

ABSTRACT

The free access to large-scale public databases, together with the fast progress of deep learning techniques, in particular Generative Adversarial Networks, have led to the generation of very realistic fake content with its corresponding implications towards society in this era of fake news.

This survey provides a thorough review of techniques for manipulating face images including DeepFake methods, and methods to detect such manipulations. In particular, four types of facial manipulation are reviewed: *i*) entire face synthesis, *ii*) identity swap (DeepFakes), *iii*) attribute manipulation, and *iv*) expression swap. For each manipulation group, we provide details regarding manipulation techniques, existing public databases, and key benchmarks for technology evaluation of fake detection methods, including a summary of results from those evaluations. Among all the aspects discussed in the survey, we pay special attention to the latest generation of DeepFakes, highlighting its improvements and challenges for fake detection.

In addition to the survey information, we also discuss open issues and future trends that should be considered to advance in the field.

1. Introduction

Fake images and videos including facial information generated by digital manipulation, in particular with DeepFake methods [1], have become a great public concern recently [2,3]. The very popular term “DeepFake” is referred to a deep learning based technique able to create fake videos by swapping the face of a person by the face of another person. This term was originated after a Reddit user named “deepfakes” claimed in late 2017 to have developed a machine learning algorithm that helped him to transpose celebrity faces into porn videos [4]. In addition to fake pornography, some of the more harmful usages of such fake content include fake news, hoaxes, and financial fraud. As a result, the area of research traditionally dedicated to general media forensics [5–11], is being invigorated and is now dedicating growing efforts for detecting facial manipulation in image and video [12]. Part of these renewed efforts in fake face detection are built around past research in biometric anti-spoofing [13–15] and modern data-driven deep learning [16,17]. The growing interest in fake face detection is demonstrated through the increasing number of workshops in top conferences [18–22], international projects such as MediFor funded by the Defense Advanced Research Project Agency (DARPA), and competitions such as the recent Media Forensics Challenge (MFC2018)¹ and the Deepfake Detection Challenge (DFDC)² launched by the National Institute of Standards and Technology (NIST) and Facebook, respectively.

Traditionally, the number and realism of facial manipulations have been limited by the lack of sophisticated editing tools, the domain expertise required, and the complex and time-consuming process involved. For example, an early work in this topic [23] was able to modify the lip motion of a person speaking using a different audio track, by making connections between the sounds of the audio track and the shape of the subject’s face. However, from these early works up to date, many things have rapidly evolved in the last years. Nowadays, it is becoming increasingly easy to automatically synthesise non-existent faces or manipulate a real face of one person in an image/video, thanks to: *i*) the accessibility to large-scale public data, and *ii*) the evolution of deep learning techniques that eliminate many manual editing steps such as Autoencoders (AE) and Generative Adversarial Networks (GAN) [24,25]. As a result, open software and mobile application such as ZAO³ and FaceApp⁴ have been released opening the door to anyone to create fake images and videos, without any experience in the field needed.

In response to those increasingly sophisticated and realistic manipulated content, large efforts are being carried out by the research community to design improved methods for face manipulation detection. Traditional fake detection methods in media forensics have been com-

* Corresponding author.

E-mail address: ruben.tolosana@uam.es (R. Tolosana).

¹ <https://www.nist.gov/itl/iad/mig/media-forensics-challenge-2018>.

² <https://deepfakedetectionchallenge.ai/>.

³ <https://apps.apple.com/cn/app/id1465199127>.

⁴ <https://apps.apple.com/gb/app/faceapp-ai-face-editor/id1180884341>.

monly based on: *i*) in-camera fingerprints, the analysis of the intrinsic fingerprints introduced by the camera device, both hardware and software, such as the optical lens [26], colour filter array and interpolation [27,28], and compression [29,30], among others, and *ii*) out-camera fingerprints, the analysis of the external fingerprints introduced by editing software, such as copy-paste or copy-move different elements of the image [31,32], reduce the frame rate in a video [33–35], etc. However, most of the features considered in traditional fake detection methods are highly dependent on the specific training scenario, being therefore not robust against unseen conditions [6,8,16]. This is of special importance in the era we live in as most media fake content is usually shared on social networks, whose platforms automatically modify the original image/video, e.g., through compression and resize operations [12].

This survey provides an in-depth review of digital manipulation techniques applied to facial content due to the large number of possible harmful applications, e.g., the generation of fake news that would provide misinformation in political elections and security threats [36,37]. Specifically, we cover four types of manipulations: *i*) entire face synthesis, *ii*) identity swap, *iii*) attribute manipulation, and *iv*) expression swap. These four main types of face manipulation are well established by the research community, receiving most attention in the last few years. Besides, we also review in this survey some other challenging and dangerous face manipulation that are not so popular yet like face morphing.

Finally, for completeness, we would like to highlight other recent surveys in the field. In [38], the authors cover the topic of DeepFakes from a general perspective, proposing the R.E.A.L framework to manage DeepFake risks. In addition, Verdoliva has recently surveyed in [39] traditional manipulation and fake detection approaches considered in general media forensics, and also the latest deep learning techniques. The present survey complements [38] and [39] with a more detailed review of each facial manipulation group, including manipulation techniques, existing public databases, and key benchmarks for technology evaluation of fake detection methods, including a summary of results from those evaluations. In addition, we pay special attention to the latest generation of DeepFakes, highlighting its improvements and challenges for fake detection.

The remainder of the article is organised as follows. We first provide in Section 2 a general description of different types of facial manipulation. Then, from Section 3 to Section 6 we describe the key aspects of each type of facial manipulation including public databases for research, detection methods, and benchmark results. Section 7 focuses on other interesting types of face manipulation techniques not covered in previous sections. Finally, we provide in Section 8 our concluding remarks, highlighting open issues and future trends.

2. Types of facial manipulations

Facial manipulations can be categorised in four main different groups regarding the level of manipulation. Fig. 1 graphically summarises each facial manipulation group. A description of each of them is provided below, from higher to lower level of manipulation:

- **Entire Face Synthesis:** this manipulation creates entire non-existent face images, usually through powerful GAN, e.g., through the recent StyleGAN approach proposed in [41]. These techniques achieve astonishing results, generating high-quality facial images with a high level of realism. Fig. 1 shows some examples for entire face synthesis generated using StyleGAN⁵. This manipulation could benefit many different sectors such as the video game and 3D-modelling industries, but it could also be used for harmful applications such as the creation of very realistic fake profiles in social networks in order to generate misinformation.
- **Identity Swap:** this manipulation consists of replacing the face of one person in a video with the face of another person. Two different

approaches are usually considered: *i*) classical computer graphics-based techniques such as FaceSwap⁶, and *ii*) novel deep learning techniques known as DeepFakes⁷, e.g., the recent ZAO mobile application. Very realistic videos of this type of manipulation can be seen on Youtube⁸ and obtained through commercial websites⁹. This type of manipulation could benefit many different sectors, in particular the film industry. However, in the other side, it could also be used for bad purposes such as the creation of celebrity pornographic videos, hoaxes, and financial fraud, among many others.

- **Attribute Manipulation:** this manipulation, also known as face editing or face retouching, consists of modifying some attributes of the face such as the colour of the hair or the skin, the gender, the age, adding glasses, etc [42]. This manipulation process is usually carried out through GAN such as the StarGAN approach proposed in [43]. One example of this type of manipulation is the popular FaceApp mobile application. Consumers could use this technology to try on a broad range of products such as cosmetics and makeup, glasses, or hairstyles in a virtual environment.
- **Expression Swap:** this manipulation, also known as face reenactment, consists of modifying the facial expression of the person. Although different manipulation techniques are proposed in the literature, e.g., at image level through popular GAN architectures [44], in this group we focus on the most popular techniques Face2Face and NeuralTextures [45,46], which replaces the facial expression of one person in a video with the facial expression of another person. This type of manipulation could be used with serious consequences, e.g., the popular video of Mark Zuckerberg saying things he never said¹⁰.

3. Entire face synthesis

3.1. Manipulation techniques and public databases

This manipulation creates entire non-existent face images. Table 1 summarises the main publicly available databases for research on detection of image manipulation techniques relying on entire face synthesis. Four different databases of fake images are of relevance here, all of them based on the same GAN architectures: ProGAN [48] and StyleGAN [41]. It is interesting to remark that each fake image may be characterised by a specific GAN fingerprint just like natural images are identified by a device-based fingerprint (i.e., PRNU). In fact, these fingerprints seem to be dependent not only of the GAN architecture, but also of the different instances of it [49–51].

In addition, as indicated in Table 1, it is important to note that the four mentioned databases only contain fake images generated using the GAN architectures discussed. In order to perform fake detection experiments on this manipulation group, researchers need to obtain real face images from other public databases such as CelebA [52], FFHQ [41], CASIA-WebFace [53], and VGGFace2 [54], among others.

We provide next a description of each public database. In [41], Karras *et al.* released a set of 100,000 synthetic face images, named 100K-Generated-Images¹¹. This database was generated using their proposed StyleGAN architecture, which was trained using the FFHQ dataset [41]. StyleGAN is an improved version of their previous popular approach ProGAN, which introduced a new training methodology based on improving both generator and discriminator progressively. StyleGAN proposes an alternative generator architecture that leads to an automatically learned, unsupervised separation of high-level attributes (e.g., pose and identity when trained on human faces) and stochastic variation in

⁶ <https://github.com/MarekKowalski/FaceSwap>.

⁷ <https://github.com/deepfakes/faceswap>.

⁸ <https://www.youtube.com/watch?v=Ulv0EW715rs>.

⁹ <https://deepfakesweb.com/>.

¹⁰ <https://www.bbc.com/news/technology-48607673>.

¹¹ <https://github.com/NVlabs/stylegan>.

⁵ <https://thispersondoesnotexist.com>.

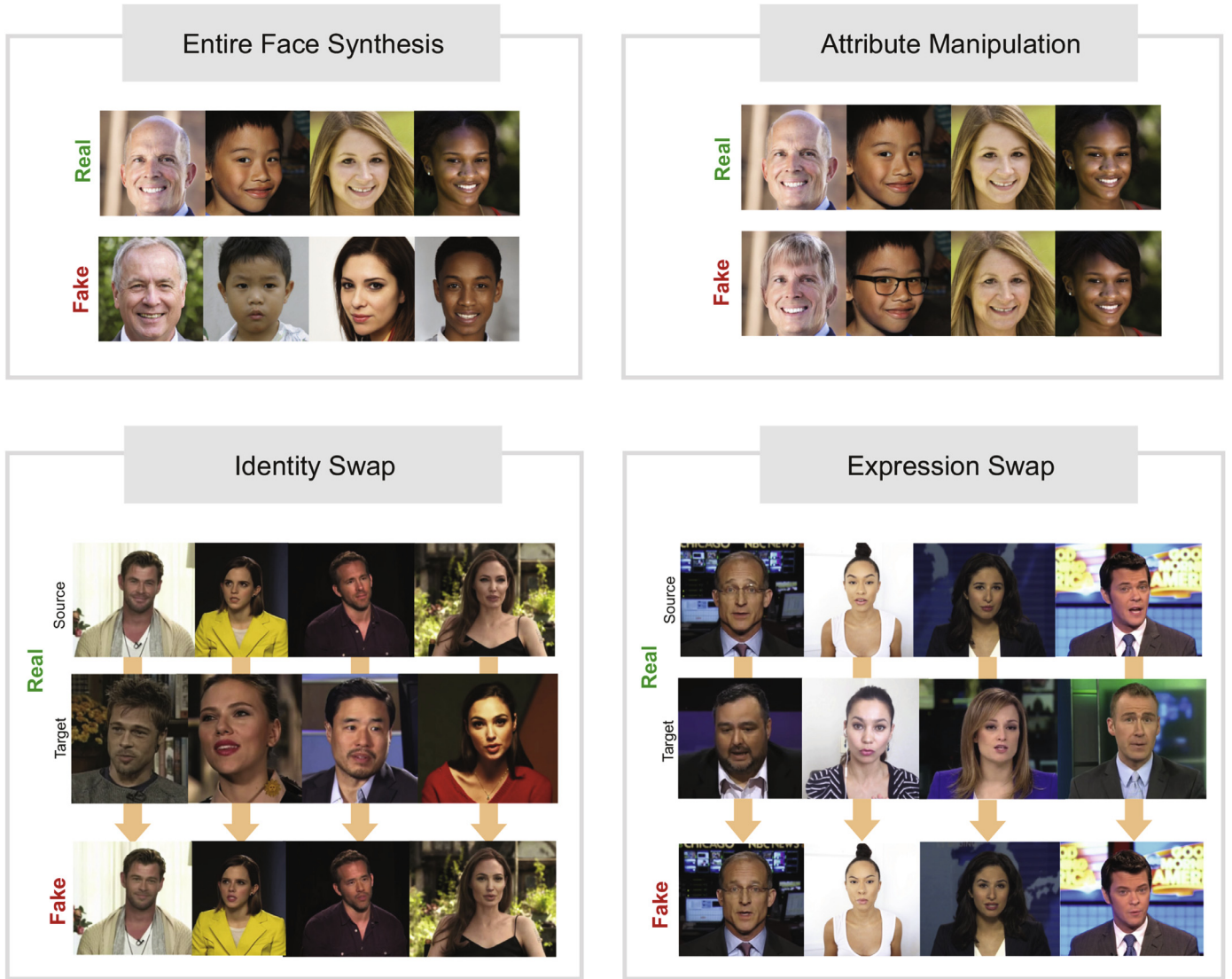


Fig. 1. Real and fake examples of each facial manipulation group. For *Entire Face Synthesis*, real images are extracted from <http://www.whichfacesreal.com/> and fake images from <https://thispersondoesnotexist.com>. For *Identity Swap*, face images are extracted from Celeb-DF database [40]. For *Attribute Manipulation*, real images are extracted from <http://www.whichfacesreal.com/> and fake images are generated using FaceApp. Finally, for *Expression Swap*, images are extracted from FaceForensics++ [12]

Table 1
Entire Face Synthesis: Publicly available databases.

Database	Real Images	Fake Images
100K-Generated-Images (2019) [41]	-	100,000 (StyleGAN)
100K-Faces (2019) [47]	-	100,000 (StyleGAN)
DFFD (2020) [17]	-	100,000 (StyleGAN) 200,000 (ProGAN)
iFakeFaceDB (2020) [16]	-	250,000 (StyleGAN) 80,000 (ProGAN)

the generated images (e.g., freckles, hair), and it enables intuitive, scale-specific control of the synthesis.

Another public database is 100K-Faces [47]. This database contains 100,000 synthetic images generated using StyleGAN. In this database, contrary to the 100K-Generated-Images database, the StyleGAN network was trained using around 29,000 photos from 69 different models, considering face images from a more controlled scenario (e.g., with a flat background). Thus, no strange artifacts created by the StyleGAN are included in the background of the images.

Recently, Dang *et al.* introduced in [17] a new database named Diverse Fake Face Dataset (DFFD). Regarding the entire face synthesis

manipulation, the authors created 100,000 and 200,000 fake images through the pre-trained ProGAN and StyleGAN models, respectively.

Finally, Neves *et al.* presented in [16] the iFakeFaceDB database. This database comprises 250,000 and 80,000 synthetic face images created with StyleGAN and ProGAN, respectively. As an additional feature in comparison to previous databases, and in order to hinder fake detectors, in this database the fingerprints produced by the GAN architectures were removed through an approach named GANprintR (GAN fingerprint Removal), while keeping very realistic appearance. Fig. 2 shows an example of a fake image directly generated with StyleGAN and its improved version after removing the GAN-fingerprint information. As a

Table 2

Entire Face Synthesis: Comparison of different state-of-the-art detection approaches. The best results achieved for each public database are remarked in **bold**. Results in *italics* indicate that they were not provided in the original work. AUC = Area Under the Curve, Acc. = Accuracy, EER = Equal Error Rate.

Study	Method	Classifiers	Best Performance	Databases (Generation)
McCloskey and Albright (2018) [55]	GAN-Pipeline Features	SVM	AUC = 70.0%	NIST MFC2018
Wang <i>et al.</i> (2019) [56]	GAN-Pipeline Features	SVM	Acc. = 84.7%	Own (InterFaceGAN, StyleGAN)
Guarnera <i>et al.</i> (2020) [57]	GAN-Pipeline Features	k-NN, SVM, LDA	Acc. = 99.81%	Own (AttGAN, GDWCT, StarGAN, StyleGAN, StyleGAN2)
Nataraj <i>et al.</i> (2019) [58]	Steganalysis Features	CNN	<i>EER = 12.3%</i> [16]	100K-Faces (StyleGAN)
Yu <i>et al.</i> (2019) [59]	Deep Learning Features	CNN	Acc. = 99.5%	Own (ProGAN, SNGAN, CramerGAN, MMDGAN)
Marra <i>et al.</i> (2019) [60]	Deep Learning Features	CNN + Incremental Learning	Acc. = 99.3%	Own (CycleGAN, ProGAN, Glow, StarGAN, StyleGAN)
Dang <i>et al.</i> (2020) [17]	Deep Learning Features	CNN + Attention Mechanism	AUC = 100% EER = 0.1%	DFFD (ProGAN, StyleGAN)
Neves <i>et al.</i> (2020) [16]	Deep Learning Features	CNN	EER = 0.3% EER = 4.5%	100K-Faces (StyleGAN) iFakeFaceDB
Hulzebosch <i>et al.</i> (2020) [61]	Deep Learning Features	CNN, AE	Acc. = 99.8%	Own (StarGAN, Glow, ProGAN, StyleGAN)



(a) Fake (b) Fake after GANprintR

Fig. 2. Examples of a fake image created using StyleGAN and its improved version after removing the GAN-fingerprint information with GANprintR [16].

result of the GANprintR step, iFakeFaceDB presents a higher challenge for advanced fake detectors compared with the other databases.

3.2. Manipulation detection

Different studies have recently evaluated the difficulty of detecting whether faces are real or artificially generated. Table 2 shows a comparison of the most relevant approaches in this area. For each study, we include information related to the method, classifiers, best performance, and databases considered. We highlight in **bold** the best results achieved for each public database. It is important to remark that in some cases, different evaluation metrics are considered, e.g., Area Under the Curve (AUC) or Equal Error Rate (EER), which complicates the comparison among the studies.

Some authors propose to analyse the internal GAN pipeline in order to detect different artifacts between real and fake images. In [55], the authors hypothesised that the colour is markedly different between real camera images and fake synthesis images. They proposed a detection system based on colour features and a linear Support Vector Machine (SVM) for the final classification, achieving a final 70.0% AUC for the best performance when evaluating with the NIST MFC2018 dataset [62].

Another interesting approach in this line was proposed in [56]. Wang *et al.* conjectured that monitoring neuron behavior could also serve as an asset in detecting fake faces since layer-by-layer neuron activation patterns may capture more subtle features that are important for the facial manipulation detection system. Their proposed approach, named FakeSpoter, extracted as features neuron coverage behaviors of real and fake faces from deep face recognition systems (i.e., VGG-Face [63], OpenFace [64], and FaceNet [65]), and then trained a SVM for the fi-

nal classification. The authors tested their proposed approach using real faces from CelebA-HQ [48] and FFHQ [41] databases and synthetic faces created through InterFaceGAN [66] and StyleGAN [41], achieving for the best performance a final 84.7% fake detection accuracy using the FaceNet model.

Better results have been recently reported in [57]. The authors proposed a fake detection system based on the analysis of the convolutional traces. Features were extracted using the Expectation Maximization algorithm [67]. Popular classifiers such as *k*-Nearest Neighbours (*k*-NN), SVM, and Linear Discriminant Analysis (LDA) were used for the final detection. Their proposed approach was tested using fake images generated through AttGAN [68], GDWCT [69], StarGAN [43], StyleGAN, and StyleGAN2 [70], achieving a final 99.81% Acc. for the best performance.

Fake detection systems inspired in steganalysis have also been studied. Nataraj *et al.* proposed in [58] a detection system based on a combination of pixel co-occurrence matrices and Convolutional Neural Networks (CNN). Their proposed approach was initially tested through a database of various objects and scenes created through CycleGAN [71]. Besides, the authors performed an interesting analysis to see the robustness of the proposed approach against fake images created through different GAN architectures (CycleGAN vs. StarGAN), with good generalisation results. This detection approach was implemented later on in [16] considering images from the 100K-Faces database, achieving an EER of 12.3% for the best fake detection performance. This result is remarked in *italics* in Table 2 to indicate that it was not provided in the original paper.

Many studies have also focused on the detection of the special fingerprints inserted by GAN architectures using pure deep learning methods. Yu *et al.* proposed in [59] an attribution network architecture to map an input image to its corresponding fingerprint image. Therefore, they learned a model fingerprint for each source (each GAN instance plus the real world), such that the correlation index between one image fingerprint and each model fingerprint serves as softmax logit for classification. Their proposed approach was tested using real faces from CelebA database [52] and synthetic faces created through different GAN approaches (ProGAN [48], SNGAN [72], CramerGAN [73], and MMDGAN [74]), achieving a final 99.5% fake detection accuracy for the best performance. However, this approach seemed not to be very robust against unseen simple image perturbation attacks such as noise, blur, cropping or compression, unless the models were re-trained again.

Related to the unseen conditions just commented, Marra *et al.* performed in [60] an interesting study in order to detect unseen types of fake generated data. Concretely, they proposed a multi-task incre-

mental learning detection method in order to detect and classify new types of GAN generated images, without worsening the performance on the previous ones. Two different solutions regarding the position of the classifier were proposed based on the successful algorithm iCaRL for incremental learning [75]: *i*) Multi-Task MultiClassifier (MT-MC), and *ii*) Multi-Task Single Classifier (MT-SC). Regarding the experimental framework, five different GAN approaches were considered in the study, CycleGAN [71], ProGAN [48], Glow [76], StarGAN [43], and StyleGAN [41]. Their proposed detection approach, based on the XceptionNet model, achieved promising results being able to correctly detect new GAN generated images.

Attention mechanisms have also been applied to further improve the training process of the detection systems. Dang *et al.* carried out in [17] a complete analysis of different types of facial manipulations. They proposed to use attention mechanisms and popular CNN models such as XceptionNet and VGG16. For the entire face synthesis manipulation, the authors achieved a final 100% AUC and around 0.1% EER considering real faces from CelebA [52], FFHQ [41], and FaceForensics++ [12] databases and fake images created through ProGAN [48] and StyleGAN [41] approaches. The impressive results achieved show the importance of novel attention mechanisms [77].

Neves *et al.* performed in [16] an in-depth experimental assessment of this type of facial manipulation considering different state-of-the-art detection systems and experimental conditions, i.e., controlled and in-the-wild scenarios. Four different fake databases were considered: *i*) 150,000 fake faces collected online¹² and based on StyleGAN architecture, *ii*) the 100K-faces public database, *iii*) 80,000 synthetic faces generated using ProGAN, and *iv*) the iFakeFaceDB database, an improved version of previous fake databases in which the GAN-fingerprint information has been removed using the GANprintR approach. In controlled scenarios, they achieved similar results as the best previous studies (EER = 0.02%). However, in more challenging scenarios in which images (real and fake) come from different sources (mismatch of datasets), a high degradation of the fake detection performance is observed. Finally, the results achieved over their public iFakeFaceDB database with an EER = 4.5% for the best fake detectors remark how challenging is iFakeFaceDB even for the most advanced manipulation detection methods. Related to this enhanced fake content, Cozzolino *et al.* proposed in [78] a similar approach based on GAN to inject camera traces into synthetic images to spoof state-of-the-art fake detectors.

Similar to [16], Hulzebosch *et al.* have recently performed in [61] an in-depth analysis of this face manipulation considering different scenarios such as cross-model, cross-data, and post-processing. Fake detectors were based on the popular Xception network and ForensicTransfer [79], which is an Autoencoder approach. In general, bad generalisation results were obtained under unseen scenarios, similar to [16].

Finally, we also include for completeness some important references to other recent studies focused on the detection of general GAN-based image manipulations, not facial ones. In particular, we refer the reader to [80,81].

4. Identity swap

4.1. Manipulation techniques and public databases

This is one of the most popular face manipulation research lines nowadays due to the great public concerns around DeepFakes [2,3]. It consists of replacing the face of one person in a video with the face of another person. Unlike the entire face synthesis manipulation, where manipulations are carried out at image level, in identity swap the goal is to generate realistic fake videos.

Since publicly available fake databases such as the UADFV database [82], up to the recent Celeb-DF and DFDC databases [40,83], many vi-

sual improvements have been carried out, increasing the realism of fake videos. As a result, identity swap databases can be divided into two different generations. Table 3 summarises the main details of each public database, grouped in each generation. As can be seen, in this type of facial manipulation both real and fake videos are usually included in the databases.

In this section, we first provide the main details of each database, to finally summarise at a higher level the key differences among the two generations.

Three different databases are grouped in the first generation. UADFV was one of the first public databases [82]. This database comprises 49 real videos from Youtube, which were used to create 49 fake videos through the FakeApp application¹³, swapping in all of them the original face with the face of Nicolas Cage. Therefore, only one identity is considered in all fake videos. Each video represents one individual, with a typical resolution of 294×500 pixels, and 11.14 seconds on average.

Korshunov and Marcel introduced in [1] the DeepfakeTIMIT database. This database comprises 620 fake videos of 32 subjects from the VidTIMIT database [85]. Fake videos were created using the public GAN-based face-swapping algorithm¹⁴. In that approach, the generative network is adopted from CycleGAN [71], using the weights of FaceNet [65]. The method Multi-Task Cascaded Convolution Networks is used for more stable detections and reliable face alignment [86]. Besides, the Kalman filter is also considered to smooth the bounding box positions over frames and eliminate jitter on the swapped face. Regarding the scenarios considered in DeepfakeTIMIT, two different qualities are considered: *i*) low quality (LQ) with images of 64×64 pixels, and *ii*) high quality (HQ) with images of 128×128 pixels. Additionally, different blending techniques were applied to the fake videos regarding the quality level.

One of the most popular databases in this type of facial manipulation is FaceForensics++ [12]. This database was introduced early 2019 as an extension of the original FaceForensics database [87], which was focused only on expression swap. FaceForensics++ contains 1000 real videos extracted from Youtube. Regarding the identity swap fake videos, they were generated using both computer graphics and DeepFake approaches (i.e., learning approach). For the computer graphics approach, the authors considered the publicly available FaceSwap algorithm¹⁵ whereas for the DeepFake approach, fake videos were created through the DeepFake FaceSwap GitHub implementation¹⁶. The FaceSwap approach consists of face alignment, Gauss Newton optimization and image blending to swap the face of the source person to the target person. The DeepFake approach, as indicated in [12], is based on two autoencoders with a shared encoder that are trained to reconstruct training images of the source and the target face, respectively. A face detector is used to crop and to align the images. To create a fake image, the trained encoder and decoder of the source face are applied to the target face. The autoencoder output is then blended with the rest of the image using Poisson image editing [88]. Regarding the figures of the FaceForensics++ database, 1000 fake videos were generated for each approach. Later on, a new dataset named DeepFakeDetection, grouped inside the 2nd generation due to its higher realism, was included in the FaceForensics++ framework with the support of Google [84]. This dataset comprises 363 real videos from 28 paid actors in 16 different scenes. Additionally, 3068 fake videos are included in the dataset based on DeepFake FaceSwap GitHub implementation. It is important to remark that for both FaceForensics++ and DeepFakeDetection databases different levels of video quality are considered, in particular: *i*) RAW (original quality), *ii*) HQ (constant rate quantization parameter equal to 23), and *iii*) LQ (constant rate quantization parameter equal to 40).

¹³ <https://www.malavida.com/en/soft/fakeapp/>.

¹⁴ <https://github.com/shaoanlu/faceswap-GAN>.

¹⁵ <https://github.com/MarekKowalski/FaceSwap>.

¹⁶ <https://github.com/deepfakes/faceswap>.

¹² <https://thispersondoesnotexist.com>.

Table 3
Identity Swap: Publicly available databases.

1st Generation		
Database	Real Videos	Fake Videos
UADFV (2018) [82]	49 (Youtube)	49 (FakeApp)
DeepfakeTIMIT (2018) [1]	-	620 (faceswap-GAN)
FaceForensics++ (2019) [12]	1,000 (Youtube)	1,000 (FaceSwap) 1,000 (DeepFake)
2nd Generation		
Database	Real Videos	Fake Videos
DeepFakeDetection (2019) [84]	363 (Actors)	3,068 (DeepFake)
Celeb-DF (2019) [40]	890 (Youtube)	5,639 (DeepFake)
DFDC Preview (2019) [83]	1,131 (Actors)	4,119 (Unknown)

This aspect simulates the video processing techniques usually applied in social networks.

Regarding the databases included in the 2nd generation, we highlight the recent Celeb-DF and DFDC databases released at the end of 2019. Li *et al.* presented in [40] the Celeb-DF database. This database aims to provide fake videos of better visual qualities, similar to the popular videos that are shared on the Internet¹⁷, in comparison to previous databases that exhibit low visual quality with many visible artifacts. Celeb-DF consists of 890 real videos extracted from Youtube, and 5639 fake videos, which were created through a refined version of a public DeepFake generation algorithm, improving aspects such as the low resolution of the synthesised faces and colour inconsistencies.

Facebook in collaboration with other companies and academic institutions such as Microsoft, Amazon, and the MIT launched at the end of 2019 a new challenge named the Deepfake Detection Challenge (DFDC) [83]. They first released a preview dataset consisting of 1131 real videos from 66 paid actors, and 4119 fake videos. Fake videos were generated using two different unknown approaches. The complete DFDC dataset was released later and comprises over 470 GB of content (real and fake)¹⁸.

Finally, to conclude this section, we discuss at a higher level the key differences among fake databases from the 1st and 2nd generations. In general, fake videos from the 1st generation are characterised by: *i*) low-quality synthesised faces, *ii*) different colour contrast among the synthesised fake mask and the skin of the original face, *iii*) visible boundaries of the fake mask, *iv*) visible facial elements from the original video, *v*) low pose variations, and *vi*) strange artifacts among sequential frames. Also, they usually consider controlled scenarios in terms of camera position and light conditions. Many of these aspects have been successfully improved in databases of the 2nd generation, not only at visual level, but also in terms of variability (in-the-wild scenarios). For example, the recent DFDC database considers different acquisition scenarios (i.e., indoors and outdoors), light conditions (i.e., day, night, etc.), distances from the person to the camera, and pose variations, among others. Fig. 3 graphically summarises the weaknesses present in identity swap databases of the 1st generation and the improvements carried out in the 2nd generation. Finally, it is also interesting to remark the larger number of fake videos included in the databases of the 2nd generation.

4.2. Manipulation detection

The development of novel methods to detect identity swap manipulations is continuously evolving. Table 4 provides a comparison of the most relevant detection approaches in this area. For each study we include information related to the method, classifiers, best performance, and databases for research. We highlight in **bold** the best results achieved for each public database. It is important to remark that in some cases, different evaluation metrics are considered (e.g., AUC and

EER), which complicates the comparison among studies. Finally, the results highlighted in *italics* indicate the generalisation capacity of the detection systems against different unseen databases, i.e., those databases were not considered for training. These results have been extracted from [40] and were not included in the original publications.

The first studies in this area focused on the audio-visual artifacts existed in the 1st generation of fake videos. Korshunov and Marcel evaluated in [1] baseline approaches based on the inconsistencies between lip movements and audio speech, as well as several variations of image-based systems often used in biometrics. For the first case, they considered Mel-Frequency Cepstral Coefficients (MFCCs) as audio features and distances between mouth landmarks as visual features. Principal Component Analysis (PCA) was then used to reduce the dimensionality of the blocks of features, and finally Recurrent Neural Networks (RNNs) based on Long Short-Term Memory (LSTM) to detect real of fake videos (based on [101]). For the second case, they evaluated detection approaches based on: *i*) raw faces as features, and *ii*) image quality measures (IQM) [102]. In particular, they used a set of 129 features related to measures like signal to noise ratio, specularly, blurriness, etc. PCA with LDA, or SVM were considered for the final classification. Their proposed detection approach based on IQM + SVM provided the best results, with a final 3.3% and 8.9% EER for the LQ and HQ scenarios of the DeepfakeTIMIT database, respectively.

In this line, Matern *et al.* proposed in [89] fake detection systems based on relatively simple visual aspects such as eye colour, missing reflections, and missing details in the eye and teeth areas. Two different classifiers were considered in this analysis: *i*) a logistic regression model, and *ii*) a Multilayer Perceptron (MLP) [103]. Their proposed approach was tested using a private database, achieving a final 85.1% AUC for the MLP system.

Fake detection systems based on facial expressions and head movements have also been proposed in the literature. Yang *et al.* observes in [90] that some DeepFakes are created by splicing synthesised face regions into the original image, and in doing so, introducing errors that can be revealed when 3D head poses are estimated from the face images. Thus, they performed an study based on the differences between head poses estimated using a full set of facial landmarks (68 extracted from DLib [104]) and those in the central face regions to differentiate DeepFakes from real videos. Once these features are extracted and normalised (mean and standard deviation), a SVM is considered for the final classification. Their proposed approach was originally evaluated with the UADFV database, achieving a final 89.0% AUC. However, this pre-trained model (using UADFV database) seems not to generalise very well to other databases as depicted in Table 4.

Another interesting approach in this line was proposed by Agarwal and Farid in [91]. They proposed a detection system based on both facial expressions and head movements. For the feature extraction, the OpenFace2 toolkit was considered [105], obtaining an intensity and occurrence for 18 different facial action units related to movements of facial muscles such as cheek raiser, nose wrinkle, mouth stretch, etc. Additionally, four features related to head movements were considered.

¹⁷ https://www.youtube.com/channel/UCKpH0CKltc73e4wh0_pgL3g.

¹⁸ <https://www.kaggle.com/c/deepfake-detection-challenge>.



Fig. 3. Graphical representation of the weaknesses present in identity swap databases of the 1st generation and the improvements carried out in the 2nd generation, not only at visual level, but also in terms of variability (in-the-wild scenarios). Fake images are extracted from: UADFV and FaceForensics++ (1st generation) [12,82]; Celeb-DF and DFDC (2nd generation) [40,83].

As a result, each 10-second video clip is reduced to a feature vector of dimension 190 using the Pearson correlation to measure the linearity between features. Finally, the authors considered a SVM for the final classification. Regarding the experimental framework, the authors built their own database based on videos downloaded from YouTube of persons of interest talking in a formal setting, for example, weekly address, news interview, and public speech. In most videos the person is primarily facing towards the camera. Regarding the DeepFake videos, the

authors trained one GAN per person based on faceswap-GAN¹⁹. Their proposed approach achieved a final 96.3% AUC as the best fake detection performance, being robust against new contexts and manipulation techniques.

Eye blinking [106] has also been studied to detect fake videos. In [92], the authors proposed an algorithm called DeepVision to analyse

¹⁹ <https://github.com/shaoanlu/faceswap-GAN>.

Table 4

Identity Swap: Comparison of different state-of-the-art detection approaches. The best results achieved for each public database are remarked in **bold**. Results in *italics* indicate that they were published in [40], but not in the original work. FF++ = FaceForensics++ , AUC = Area Under the Curve, Acc. = Accuracy, EER = Equal Error Rate, TCR = True Classification Rates.

Study	Method	Classifiers	Best Performance	Databases
Korshunov and Marcel (2018) [1]	Audio-Visual Features	PCA+RNN PCA + LDA, SVM	EER = 3.3% EER = 8.9%	DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ)
Matern et al. (2019) [89]	Visual Features	Logistic Regression MLP	AUC = 85.1% AUC = 70.2% AUC = 77.0% AUC = 77.3% AUC = 78.0% AUC = 66.2% AUC = 55.1%	Own UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ / DFD DFDC Preview Celeb-DF
Yang et al. (2019) [90]	Head Pose Features	SVM	AUC = 89.0% AUC = 55.1% AUC = 53.2% AUC = 47.3% AUC = 55.9% AUC = 54.6%	UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ / DFD DFDC Preview Celeb-DF
Agarwal and Farid (2019) [91]	Head Pose and Facial Features	SVM	AUC = 96.3%	Own (FaceSwap, HQ)
Jung et al. (2020) [92]	Eye Blinking	Distance	Acc. = 87.5%	Own
Li et al. (2019) [40,93]	Face Warping Features	CNN	AUC = 97.7% AUC = 99.9% AUC = 99.7% AUC = 93.0% AUC = 75.5% AUC = 64.6%	UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ / DFD DFDC Preview Celeb-DF
Afchar et al. (2018) [94]	Mesoscopic Features	CNN	Acc. = 98.4% AUC = 84.3% AUC = 87.8% AUC = 68.4% Acc. ≈ 90.0% Acc. ≈ 94.0% Acc. ≈ 98.0% Acc. ≈ 83.0% Acc. ≈ 93.0% Acc. ≈ 96.0% AUC = 75.3% AUC = 54.8%	Own UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ (DeepFake, LQ) FF++ (DeepFake, HQ) FF++ (DeepFake, RAW) FF++ (FaceSwap, LQ) FF++ (FaceSwap, HQ) FF++ (FaceSwap, RAW) DFDC Preview Celeb-DF
Zhou et al. (2018) [95]	Steganalysis Features + Deep Learning Features	CNN SVM	AUC = 85.1% AUC = 83.5% AUC = 73.5% AUC = 70.1% AUC = 61.4% AUC = 53.8%	UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ / DFD DFDC Preview Celeb-DF
Rössler et al. (2019) [12]	Mesoscopic Features Steganalysis Features Deep Learning Features	CNN	Acc. ≈ 94.0% Acc. ≈ 98.0% Acc. ≈ 100.0% Acc. ≈ 93.0% Acc. ≈ 97.0% Acc. ≈ 99.0%	FF++ (DeepFake, LQ) FF++ (DeepFake, HQ) FF++ (DeepFake, RAW) FF++ (FaceSwap, LQ) FF++ (FaceSwap, HQ) FF++ (FaceSwap, RAW)
Nguyen et al. (2019) [96]	Deep Learning Features	AE + Multi-Task Learning	AUC = 65.8% AUC = 62.2% AUC = 55.3% AUC = 76.3% EER = 15.1% AUC = 53.6% AUC = 54.3%	UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ / DFD FF++ (FaceSwap, HQ) DFDC Preview Celeb-DF
Nguyen et al. (2019) [97]	Deep Learning Features	Capsule Networks	AUC = 61.3% AUC = 78.4% AUC = 74.4% AUC = 96.6% AUC = 53.3% AUC = 57.5%	UADFV DeepfakeTIMIT (LQ) DeepfakeTIMIT (HQ) FF++ / DFD DFDC Preview Celeb-DF
Dang et al. (2019) [17]	Deep Learning Features	CNN + Attention Mechanism	AUC = 99.4% EER = 3.1%	DFFD
Dolhansky et al. (2019) [83]	Deep Learning Features	CNN	Precision = 93.0% Recall = 8.4%	DFDC Preview
Wang and Dantcheva (2020) [98]	Deep Learning Features	3DCNN	TCR = 95.13% TCR = 92.25%	FF++ (DeepFake, LQ) FF++ (FaceSwap, LQ)
Güera and Delp (2018) [99]	Image + Temporal Features	CNN + RNN	Acc. = 97.1%	Own
Sabir et al. (2019) [98]	Image + Temporal Features	CNN + RNN	AUC = 96.9% AUC = 96.3%	FF++ (DeepFake, LQ) FF++ (FaceSwap, LQ)
Tolosana et al. (2020) [100]	Facial Regions Features	CNN	AUC = 100.0% AUC = 99.4% AUC = 91.0% AUC = 83.6%	UADFV FF++ (FaceSwap, HQ) DFDC Preview Celeb-DF

changes in the blinking patterns. Their approach was based on the fusion of Fast-HyperFace [107] and Eye-Aspect-Ratio (EAR) [108] to detect the face and obtain the eye aspect ratio. Finally, features based on blinking count and period were extracted to decide whether the video is real or fake. This approach achieved a final 87.5% accuracy over a proprietary database.

Another interesting research line is based on the detection of the artifacts included by the face manipulation pipeline. In [93], Li and Lyu hypothesised that some DeepFake algorithms can only create images of limited resolution, which need to be further warped to match the original faces in the source video. Such transforms leave distinctive artifacts in the resulting DeepFake videos. Thus, the authors proposed a detection system based on CNNs in order to detect the presence of such artifacts from the detected face regions and the surrounding areas. Four different CNN models were trained from scratch: VGG16 [109], ResNet50, ResNet101, and ResNet152 [110]. Their proposed detection approach was tested using the UADFV and DeepfakeTIMIT databases, outperforming the state of the art for those databases.

Li *et al.* proposed later on in [40] an improved version of the work presented in [93]. In this case, the authors included a new spatial pyramid pooling module to better handle the variations in the resolution [111]. This detection approach was evaluated using different databases, achieving state-of-the-art results in some of them.

Approaches based on mesoscopic and steganalysis features have also been proposed in the literature. Afchar *et al.* proposed in [94] two different networks composed of few layers in order to focus on the mesoscopic properties of the images: *i*) a CNN network comprised of 4 convolutional layers followed by a fully-connected layer (Meso-4), and *ii*) a modification of Meso-4 consisted of a variant of the Inception module introduced in [112], named MesoInception-4. Their proposed approach was originally tested against DeepFakes using a private database, achieving a 98.4% of fake detection accuracy for the best performance. That pre-trained detection model was tested against unseen databases in [40], proving to be a robust approach in some cases such as with the FaceForensics++ database.

Zhou *et al.* proposed a two-stream network for face manipulation detection. In particular, the authors considered a fusion of two streams: *i*) a face classification stream based on the CNN GoogLeNet [112] to detect whether a face image is fake or not, and *ii*) a path triplet stream that is trained using steganalysis features of images patches with a triplet loss, and a SVM for the classification. The initial system was trained to detect expression swap manipulations. Nevertheless, Li *et al.* evaluated in [40] the generalisation capacity of the pre-trained model (trained using SwapMe app) to detect identity swap manipulations, resulting to be one of the most robust approaches against the recent Celeb-DF database [40].

An exhaustive analysis of different fake detection methods was carried out by Rössler *et al.* using FaceForensics++ database [12]. Five different detection systems were evaluated: *i*) a CNN-based system trained through handcrafted steganalysis features [113], *ii*) a CNN-based system whose convolution layers are specifically designed to suppress the high-level content of the image [114], *iii*) a CNN-based system with a global pooling layer that computes four statistics (mean, variance, maximum, and minimum) [115], *iv*) the CNN MesoInception-4 detection system described in [94], and finally *v*) the CNN-based system XceptionNet [116] pre-trained using ImageNet database [117] and re-trained for the face manipulation detection task. In general, the detection system based on XceptionNet architecture provided the best results in both types of manipulation methods, DeepFakes and FaceSwap. In addition, the detection systems were evaluated considering different video quality levels in order to simulate the video processing of many social networks. In this real scenario, the accuracy of all detection systems decreased when lowering the video quality, remarking how challenging is this task in real scenarios.

Recent deep learning methods considered in computer vision have been applied to further improve the detection of identity swap manipulations.

In [96], Nguyen *et al.* proposed a CNN system that uses multi-task learning to simultaneously detect fake videos and locate the manipulated regions. They considered a detection system based on an auto-encoder. Concretely, they proposed to use a Y-shaped decoder in order to share valuable information between the classification, segmentation, and reconstruction tasks, improving the overall performance by reducing the loss. Their proposed approach was evaluated with the FaceSwap manipulation method for the FaceForensics++ database [87], achieving a best performance of 15.07% EER, far from other detection approaches. In addition, this model seems not to generalise very well for other databases, with results below 80% AUC.

Later on, the same authors presented in [97] a new fake detection system based on the recent Capsule Networks. This approach uses fewer parameters than traditional CNN with similar performance [118–120]. The proposed detection system was originally evaluated using FaceForensics++ database with accuracies higher than 90%. The same pre-trained detection model was tested against unseen databases in [40], showing poor generalisation results, as it happens in most fake detection systems.

Attention mechanisms have also been applied to further improve the training process of the detection systems. Dang *et al.* performed in [17] a thorough analysis of different face manipulations. They proposed a detection system based on CNN and attention mechanisms to process and improve the feature maps of the classifier model. Their proposed attention map can be implemented easily and inserted into existing backbone networks, through the inclusion of a single convolution layer, its associated loss functions, and masking the subsequent high-dimensional features. Their proposed detection approach was tested with the DFFD database (based on a combination of the previous FaceForensics++ databases and a collection of videos from the Internet). In particular, for identity swap detection, their proposed approach achieved an AUC of 99.43% and EER of 3.1%. Despite of the fact that it is difficult to provide a fair comparison among studies as different experimental protocols are considered, it is clear that their detection approach provides state-of-the-art results.

In [83], in addition to the description of the DFDC database, the authors provided baseline results using three simple detection systems: *i*) a small CNN model composed of 6 convolution layers and 1 fully-connected layer to detect low-level image manipulations, *ii*) an XceptionNet model trained using only face images, and *iii*) an XceptionNet model trained using the full image. The detection system based on XceptionNet, considering only the face image (not the full image), provided the best results with 93.0% precision and 8.4% recall.

Deep learning approaches based on 3DCNN were studied in [121] in order to consider both spatial and motion information. In particular, the authors proposed fake detectors based on I3D [122] and 3D ResNet [123] approaches, achieving promising results on the low quality videos of the FaceForensics++ database.

Detection systems based not only on features at image level, but also at temporal level, along the frames of the video, have also been studied in the literature. Güera and Delp proposed in [99] a temporal-aware pipeline to automatically detect fake videos. They considered a combination of CNNs and RNNs. For the CNN, the authors used InceptionV3 [124] pre-trained using ImageNet database [117]. For the RNN system, they considered a LSTM model composed of one hidden layer with 2048 memory blocks. Finally, two fully-connected layers were included, providing the probabilities of the frame sequence being either real or fake. Their proposed approach was evaluated using a proprietary database with a final 97.1% accuracy.

In this line, Sabir *et al.* proposed a method to detect fake videos based on using the temporal information present in the stream [98]. The intuition behind this model is to exploit temporal discrepancies across frames. Thus, they considered a recurrent convolutional network similar to [99], trained in this study end-to-end instead of using a pre-trained model. Their proposed detection approach was tested through FaceForensics++ database, achieving AUC results of 96.9% and 96.3%

for the DeepFake and FaceSwap methods, respectively. Only the low-quality videos were considered in the analysis.

Finally, the discriminative power of each facial region for the detection of fake videos was studied in [100]. The authors considered a fake detection system based on XceptionNet. Databases from both 1st and 2nd generations were considered in the experimental framework, concluding that poor fake detection results are achieved in the latest DeepFake video databases of the 2nd generation compared with the 1st generation, with results of 91.0% and 83.6% AUC for the DFDC Preview and Celeb-DF databases, respectively. It is important to highlight that, contrary to [40], a separate fake detection system was specifically trained for each database.

In conclusion, although many different approaches have been proposed in the literature, they all show poor generalisation results to unseen databases, as indicated in Table 4. In addition, we also highlight the poor detection results achieved by most approaches on the DeepFake databases of the 2nd generation with results below 60% AUC.

5. Attribute manipulation

5.1. Manipulation techniques and public databases

This face manipulation consists of modifying in an image some attributes of the face such as the colour of the hair or the skin, the gender, the age, adding glasses, etc. Despite the success of GAN-based frameworks for general image translations and manipulations [43,71,125–129], and in particular for face attribute manipulations [43,44,68,130–134], few databases are publicly available for research in this area, to the best of our knowledge. The main reason is that the code of most GAN approaches are publicly available, so researchers can easily generate their own fake databases as they like. Therefore, this section aims to highlight the latest GAN approaches in the field, from older to closer in time, providing also the link to their corresponding codes.

In [130], the authors introduced the Invertible Conditional GAN (IcGAN)²⁰ for complex image editing as the union of an encoder used jointly with a conditional GAN (cGAN) [135]. This approach provides accurate results in terms of attribute manipulation. However, it seriously changes the face identity of the person.

Lample *et al.* proposed in [133] an encoder-decoder architecture that is trained to reconstruct images by disentangling the salient information of the image and the attribute values directly in the latent space²¹. However, as it happens with the IcGAN approach, the generated images may lack some details or present unexpected distortions.

An enhanced approach named StarGAN²² was proposed in [43]. Before the StarGAN approach, many studies had shown promising results in image-to-image translations for two domains in general. However, few studies had focused on handling more than two domains. In that case a direct approach would be to build different models independently for every pair of image domains. StarGAN proposed a novel approach able to perform image-to-image translations for multiple domains using only a single model. The authors trained a conditional attribute transfer network via attribute classification loss and cycle consistency loss. Good visual results were achieved compared with previous approaches. However, it sometimes includes undesired modifications from the input face image such as the colour of the skin.

Almost at the same time He *et al.* proposed in [68] attGAN²³, a novel approach that removes the strict attribute-independent constraint from the latent representation, and just applies the attribute-classification constraint to the generated image to guarantee the correct change of the attributes. AttGAN provides state-of-the-art results on realistic attribute manipulation with other facial details well preserved.

One of the latest approaches proposed in the literature is STGAN²⁴ [44]. In general, attribute manipulation can be tackled by incorporating an encoder-decoder or GAN. However, as commented Liu *et al.* [44], the bottleneck layer in the encoder-decoder usually provides blurry and low quality manipulation results. To improve this, the authors presented and incorporated selective transfer units with an encoder-decoder for simultaneously improving the attribute manipulation ability and the image quality. As a result, STGAN has recently outperformed the state of the art in attribute manipulation.

Despite of the fact that the code of most attribute manipulation approaches are publicly available, the lack of public databases and experimental protocols results crucial when comparing among different manipulation detection approaches. Otherwise, it is not possible to perform a fair comparison among studies. Up to now, to the best of our knowledge, the DFFD database [17] seems to be the only public database that considers this type of facial manipulations. This database comprises 18,416 and 79,960 fake images generated through FaceApp and StarGAN approaches, respectively.

5.2. Manipulation detection

Attribute manipulations have been originally studied in the field of face recognition in order to see how robust biometric systems are against physical factors such as plastic surgery, cosmetics, makeup or occlusions [136–141]. However, it has been the recent success of mobile applications such as FaceApp that has motivated the research community to detect digital face attribute manipulations. Table 5 provides a comparison of the most relevant approaches in this area. We include for each study information related to the method, classifiers, best performance, and databases for research.

Some authors propose to analyse the internal GAN pipeline to detect different artifacts between real and manipulated images. Similar to the entire face synthesis manipulations, Wang *et al.* conjectured in [56] that monitoring neuron behavior could also serve as an asset in detecting fake faces since layer-by-layer neuron activation patterns may capture more subtle features that are important for the facial manipulation detection system. Their proposed approach, named FakeSpotter, extracted as features neuron coverage behaviors of real and fake faces from deep face recognition systems (VGG-Face [63], OpenFace [64], and FaceNet [65]), and then trained a SVM for the final classification. The authors tested their proposed approach using real faces from CelebA-HQ [48] and FFHQ [41] databases and synthetic faces created through InterFaceGAN [66] and StyleGAN [41], achieving for the best performance a final 84.7% manipulation detection accuracy using the FaceNet model.

Fake detection systems inspired in steganalysis have also been studied. As described in Section 3.2 for the entire face synthesis, Nataraj *et al.* proposed in [58] a detection system based on the combination of pixel co-occurrence matrices and CNN. They created a new fake dataset based on attribute manipulations using the StarGAN approach [43] trained through the CelebA database [52], achieving a final 99.4% accuracy for the best result.

Many studies have also focused on pure deep learning methods, either feeding the networks with face patches or with the complete face. In [142], Bharati *et al.* proposed a deep learning approach based on a Restricted Boltzmann Machine (RBM) in order to detect digital retouching of face images. The input of the detection system consisted of face patches in order to learn discriminative features to classify each image as original or retouched. Regarding the databases, the authors generated two fake databases from the original ND-IIITD database (collection B [148]) and a set of celebrity facial images downloaded from the Internet. Fake images were generated using the professional software PortraitPro Studio Max²⁵, considering aspects such as skin texture, shape of

²⁰ <https://github.com/Guim3/IcGAN>.

²¹ <https://github.com/facebookresearch/FaderNetworks>.

²² <https://github.com/yunjey/stargan/blob/master/README.md>.

²³ <https://github.com/LynnHo/AttGAN-Tensorflow>.

²⁴ <https://github.com/csmlu/STGAN>.

²⁵ <https://www.anthropic.com/portraitpro/>.

Table 5

Attribute Manipulation: Comparison of different state-of-the-art detection approaches. The best results achieved for each public database are remarked in **bold**. AUC = Area Under the Curve, Acc. = Accuracy, EER = Equal Error Rate.

Study	Method	Classifiers	Best Performance	Databases (Generation)
Wang <i>et al.</i> (2019) [56]	GAN-Pipeline Features	SVM	Acc. = 84.7%	Own (InterFaceGAN/StyleGAN)
Nataraj <i>et al.</i> (2019) [58]	Steganalysis Features	CNN	Acc. = 99.4%	Own (StarGAN/CycleGAN)
Bharati <i>et al.</i> (2016) [142]	Deep Learning Features (Face Patches)	RBM	Overall Acc. = 96.2% Overall Acc. = 87.1%	Own (Celebrity Retouching, ND-IIITD Retouching)
Jain <i>et al.</i> (2019) [143]	Deep Learning Features (Face Patches)	CNN + SVM	Overall Acc. = 99.6% Overall Acc. = 99.7%	Own (ND-IIITD Retouching, StarGAN)
Tariq <i>et al.</i> (2018) [144]	Deep Learning Features	CNN	AUC = 99.9% AUC = 74.9%	Own (ProGAN, Adobe Photoshop)
Dang <i>et al.</i> (2019) [17]	Deep Learning Features	CNN + Attention Mechanism	AUC = 99.9% EER = 1.0%	DFFD (FaceApp/StarGAN)
Wang <i>et al.</i> (2019) [145]	Deep Learning Features	DRN	AP = 99.8%	Own (Adobe Photoshop)
Marra <i>et al.</i> (2019) [60]	Deep Learning Features	CNN + Incremental Learning	Acc. = 99.3%	Own (Glow/StarGAN)
Zhang <i>et al.</i> (2019) [146]	Spectrum Domain Features	GAN Discriminator	Acc. = 100%	Own (StarGAN/CycleGAN)
Rathgeb <i>et al.</i> (2020) [147]	PRNU Features	Score-Level Fusion	EER = 13.7%	Own (5 Public Apps)

eyes, nose, lips and overall face, prominence of smile, lip shape, and eye colour. Their proposed approach achieved overall accuracies for manipulation detection of 96.2% and 87.1% for the celebrity and ND-IIITD retouching databases, respectively.

A similar approach based on non-overlapping face patches was presented in [143]. Jain *et al.* proposed a CNN feature extractor composed of 6 convolutional layers and 2 fully-connected layers. Additionally, residual connections were considered inspired by a ResNet architecture [110]. Finally, a SVM was used for the final classification. Regarding the experimental framework, the ND-IIITD retouched database presented in [142] was considered. Additionally, the authors considered fake images created through the StarGAN approach [43], trained using the CelebA database [52]. In general, good detection results were achieved in both manipulation approaches, achieving almost 100% manipulation detection accuracy.

Deep learning methods based on the complete face have been further studied in the literature, achieving in general very good results. Tariq *et al.* evaluated in [144] the use of different CNN architectures such as VGG16 [63], VGG19 [63], ResNet [110], or XceptionNet [116], among others. For the real face images, the CelebA database [52] was used. Regarding the fake images, two different approaches were considered: *i*) machine approaches based on GAN, in particular ProGAN [48], and *ii*) manual approach based on Adobe Photoshop CS6, including manipulations such as makeup, glasses, sunglasses, hair, and hats. For the experimental evaluation, different sizes of the images were considered (from 32×32 to 256×256 pixels). A final 99.99% AUC was obtained for the machine-created scenario whereas for the human-created scenario this value decreased to a final 74.9% AUC for the best CNN model. Thus, a high degradation of the manipulation detection performance was observed between machine- and human-created fake images.

Attention mechanisms have also been applied to further improve the training process of the detection systems. As described in previous sections, Dang *et al.* developed in [17] a system able to detect different types of fakes. They used attention mechanisms to process and improve the feature maps of CNN models. Regarding the attribute manipulations, two different approaches were considered: *i*) fake images created through the public FaceApp software, with up to 28 different available

filters considering aspects such as hair, age, glasses, beard, and skin colour, among others; and *ii*) fake images created through the StarGAN approach [43], with up to 40 different filters. Their proposed approach was tested using their novel database DFFD, achieving very good results close to 1.0% EER (and 99.9% of AUC).

Wang *et al.* carried out in [145] an interesting research using publicly available commercial software from Adobe Photoshop (Face-Aware Liquify tool [149]) in order to synthesise new faces, and also a professional artist in order to manipulate 50 real photographs. The authors began running a human study through Amazon Mechanical Turk (AMT), showing real and fake images to the participants and asking them to classify each image into one of the classes. The results achieved remark how challenging the task is for humans, with a final 53.5% of accuracy, close to chance (50%). After the human study, the authors proposed two different automatic models: *i*) a global classification model based on Dilated Residual Networks (DRN) to predict whether the face has been warped or not, and *ii*) a local warp predictor based on the optical flow field in order to identify where manipulation occurs, and reverse them. The PWC-Net approach proposed in [150] was considered to compute the flow from original to manipulated and vice versa. Performances of 99.8% and 97.4% for automatic and manual face synthesis manipulation were achieved.

The work [60] by Marra *et al.* also described in Section 3.2 was able to correctly perform discrimination when new GANs were presented to the network and achieved a 99.3% accuracy for their proposed manipulation detection approach, based on the XceptionNet model.

A detection system based on features extracted from the spectrum domain, rather than the raw image pixels, was presented by Zhang *et al.* in [146]. Given an image as input, they applied a 2D DFT to each of the RGB channels, getting one frequency image per channel. Regarding the classifier, they proposed AutoGAN, which is a GAN simulator that can synthesise GAN artifacts in any image without needing to access any pre-trained GAN model. The generalisation capacity of their proposed approach was tested using unseen GAN models. In particular, StarGAN [43] and GauGAN [126] were considered in the evaluation. For the StarGAN approach, good detection results were achieved using the frequency domain (100%). However, for the GauGAN approach, a high

degradation of the system performance, 50% accuracy, was observed. The authors claimed that this was produced due to the generator of the GauGAN is drastically different from the CycleGAN (used in training).

Finally, Rathgeb *et al.* proposed in [147] a detection system based on Photo Response Non-Uniformity (PRNU). Specifically, scores obtained from the analysis of spatial and spectral features extracted from PRNU patterns across image cells were fused. Their proposed approach was evaluated over a private database created using 5 different mobile applications, achieving an average 13.7% EER in manipulation detection.

To summarise this section, we can see that the core of most attribute manipulation detection systems are based on deep learning technology, providing in general very good results close to 100% accuracy, as indicated in Table 5. This is mainly produced due to the GAN-fingerprint information present in fake images. However, as indicated in the entire face synthesis manipulation, recent studies have been proposed in the literature to remove such GAN fingerprints from the fake images while keeping very realistic appearance [16,78], which represent a challenge even for the most advanced manipulation detectors.

6. Expression swap

6.1. Manipulation techniques and public databases

This manipulation, also known as face reenactment, consists of modifying the facial expression of the person. We focus on the most popular techniques Face2Face and NeuralTextures, which replace the facial expression of one person in a video with the facial expression of another person (also in a video). To the best of our knowledge, the only available database for research in this area is FaceForensics++ [12], an extension of FaceForensics [87].

Initially, the FaceForensics database was focused on the Face2Face approach [45]. This is a computer graphics approach that transfers the expression of a source video to a target video while maintaining the identity of the target person. This was carried out through manual keyframe selection. Concretely, the first frames of each video were used to obtain a temporary face identity (i.e., a 3D model), and track the expression over the remaining frames. Then, fake videos were generated by transferring the source expression parameters of each frame (i.e., 76 Blendshape coefficients) to the target video. Later on, the same authors presented in FaceForensics++ a new learning approach based on NeuralTextures [46]. This is a rendering approach that uses the original video data to learn a neural texture of the target person, including a rendering network. In particular, the authors considered in their implementation a patch-based GAN-loss as used in Pix2Pix [126]. Only the facial expression corresponding to the mouth was modified. It is important to remark that all data is available on the FaceForensics++ GitHub²⁶. In total, there are 1000 real videos extracted from Youtube. Regarding the manipulated videos, 2000 fake videos are available (1,000 videos for each considered fake approach). In addition, it is important to highlight that different video quality levels are considered, in particular: *i*) RAW (original quality), *ii*) HQ (constant rate quantization parameter equal to 23), and *iii*) LQ (constant rate quantization parameter equal to 40). This aspect simulates the video processing techniques usually applied in social networks.

In addition to the Face2Face and NeuralTexture techniques considered in expression swap manipulations at video level, different approaches have been recently proposed to change the facial expression in both images and videos. A very popular approach was presented in [151]. Averbuch-Elor *et al.* proposed a technique to automatically animate a still portrait using a video of a different subject, transferring the expressiveness of the subject of the video to the target portrait. Unlike Face2Face and NeuralTexture approaches that require videos from both input and target faces, in [151] just an image of the target is needed. In

this line, a recent approach was recently presented in [152], providing very good results in both one-shot and few-shot learning.

Finally, we also highlight other popular approaches at image level. For example, mobile applications such as FaceApp²⁷ allow to easily change the level of smiling, from happier to angrier. These approaches are based on current GAN architectures. For example, Choi *et al.* showed in [43] the potential of StarGAN to change an input image to different expression levels such as angry, happy, neutral, sad, surprised, and fearful. Other recent GAN approaches that improve both the image quality of the fake images and the control editing of the parameters are InterFaceGAN [66], UGAN [153], STGAN [44], and AttGAN [68].

6.2. Manipulation detection

This section aims to provide an overview of the expression swap detectors at video level using the FaceForensics++ database, as this is the only publicly available database for research in this area, to the best of our knowledge. Manipulations at image level (not video) can be detected using the same approaches described in Section 3.2 and 5.2.

Table 6 provides a comparison of the most relevant approaches in the area of expression swap detection. For each study we include information related to the method, classifiers, best performance, and databases. We highlight in **bold** the best results achieved for the only public database, FaceForensics++. It is important to remark that in some cases, different evaluation metrics are considered (e.g., AUC and EER), which makes it difficult to perform a fair comparison among the studies.

Some of the following methods were already discussed in Sect. 4.2 for identity swap detection. Here we summarise the results achieved by them in detecting expression swap manipulations.

Preliminary studies have focused on the visual features existed in fake videos such as the eye colour, missing reflections, etc. In [89] by Matern *et al.*, the proposed approach was tested using the FaceForensics++ database, but only the Face2Face manipulation technique, achieving a final 86.6% AUC for the best performance.

Approaches based on mesoscopic and steganalysis features have also been studied in the literature. In [94], the proposed approach was tested using the Face2Face fake videos from the FaceForensics++ database [12], achieving in general good results, especially for RAW-quality videos. The same approach was later on tested in [12] against NeuralTextures fake videos, obtaining lower accuracy results compared with the Face2Face scenario.

Recent deep learning methods have also been applied with good results. In [12], the detection system based on XceptionNet provided the best results in both Face2Face and NeuralTextures manipulations, close to 100% on RAW quality. In addition, the detection systems were evaluated considering different video quality levels in order to simulate the video processing of many social networks. In this real scenario, the accuracy of all detection systems was degraded with the video quality, as it happens in identity swap manipulations.

In [96], the proposed approach based on multi-task learning was evaluated with the FaceForensics++ database. For the Face2Face method, a 7.1% EER was achieved on HQ videos whereas for the NeuralTexture method, the EER increased a bit more to a final 7.8% EER in manipulation detection.

Attention mechanisms have been recently proposed in [17] to further improve the training process. The proposed detection approach was tested using the DFFD database, which for the expression swap manipulation is based only on data from FaceForensics++ database. The proposed approach achieved an AUC = 99.4% and EER = 3.4%.

Deep learning approaches based on 3DCNN were studied in [121] in order to consider both spatial and motion information. Similar to the identity swap manipulation, the authors proposed fake detectors based

²⁶ <https://github.com/ondyari/FaceForensics>.

²⁷ <https://apps.apple.com/gb/app/faceapp-ai-face-editor/id1180884341>.

Table 6

Expression Swap: Comparison of different state-of-the-art detection approaches. The best results achieved for each public database are remarked in **bold**. FF++ = FaceForensics++ , AUC = Area Under the Curve, Acc. = Accuracy, EER = Equal Error Rate, TCR = True Classification Rate.

Study	Method	Classifiers	Best Performance	Databases (Generation)
Matern <i>et al.</i> (2019) [89]	Visual Features	Logistic Regression, MLP	AUC = 86.6%	FF++ (Face2Face, RAW)
Afchar <i>et al.</i> (2018) [94]	Mesoscopic Features	CNN	Acc. = 83.2% Acc. = 93.4% Acc. = 96.8%	FF++ (Face2Face, LQ) FF++ (Face2Face, HQ) FF++ (Face2Face, RAW)
			Acc. \approx 75% Acc. \approx 85% Acc. \approx 95%	FF++ (NeuralTextures, LQ) FF++ (NeuralTextures, HQ) FF++ (NeuralTextures, RAW)
Rössler <i>et al.</i> (2019) [12]	Mesoscopic Features Steganalysis Features Deep Learning Features	CNN	Acc. \approx 91% Acc. \approx 98% Acc. \approx 100% Acc. \approx 81% Acc. \approx 93% Acc. \approx 99%	FF++ (Face2Face, LQ) FF++ (Face2Face, HQ) FF++ (Face2Face, RAW) FF++ (NeuralTextures, LQ) FF++ (NeuralTextures, HQ) FF++ (NeuralTextures, RAW)
Nguyen <i>et al.</i> (2019) [96]	Deep Learning Features	Autoencoder	EER = 7.1% EER = 7.8%	FF++ (Face2Face, HQ) FF++ (NeuralTextures, HQ)
Dang <i>et al.</i> (2020) [17]	Deep Learning Features	CNN + Attention Mechanism	AUC = 99.4% EER = 3.4%	FF++ (Face2Face, -)
Wang and Dantcheva (2020) [98]	Deep Learning Features	3DCNN	TCR = 90.27% TCR = 80.5%	FF++ (Face2Face, LQ) FF++ (NeuralTextures, LQ)
Sabir <i>et al.</i> (2019) [98]	Image + Temporal Features	CNN + RNN	Acc. = 94.3	FF++ (Face2Face, LQ)
Amerini <i>et al.</i> (2019) [154]	Image + Temporal Features	CNN + Optical Flow	Acc. = 81.6%	FF++ (Face2Face, -)

on I3D [122] and 3D ResNet [123] approaches, achieving promising results on the low quality videos of the FaceForensics++ database.

Another interesting line is based on the analysis of both image and temporal information. In [98], the proposed approach based on recurrent convolutional networks was tested using the FaceForensics++ database, achieving AUC results of 94.3% for the Face2Face technique. Only the low-quality videos were considered in the analysis. Finally, in [154], Amerini *et al.* proposed the adoption of optical flow fields to exploit possible inter-frame dissimilarities, using the PWC-Net approach [150]. The optical flow is a vector field computed among consecutive frames to extract apparent motion in the scene. The use of this approach is motivated as fake videos should have unnatural optical flow due to the unusual movement of lips, eyes, etc. Preliminary results were obtained using both VGG16 and ResNet50 networks, obtaining an Acc. = 81.6% for the best performance in manipulation detection.

Finally, as stated previously, most of the approaches reported here for expression swap detection have also been used for identity swap detection as reviewed in Section 4.2. In general, it seems that similar features can be learnt by the fake detectors to distinguish between real and fake content, achieving good results in both types of manipulations. We highlight the potential of novel techniques such as attention mechanisms to better guide the networks during the training process, as shown in [17], achieving AUC results of 99.4% for detecting both identity swap and expression swap manipulations.

7. Other face manipulation directions

The four classes of face manipulation techniques described before are the ones that are receiving most attention in the last few years, but they do not perfectly represent all possible face manipulations. This section discusses some other challenging and dangerous approaches in face manipulation: face morphing, face de-identification, and face synthesis based on audio or text (i.e., audio-to-video and text-to-video).

7.1. Face morphing

Face morphing is a type of face manipulation that can be used to create artificial biometric face samples that resemble the biometric information of two or more individuals [155,156]. This means that the

new morphed face image would be successfully verified against facial samples of these two or more individuals creating a serious threat to face recognition systems [157,158]. In this sense, face morphing is a different type of facial manipulation compared with the four main types covered in this survey. Also, it is worth noting that face morphing is mainly focused on creating fake samples at image level, not video such as identity swap manipulations.

There has been recently a large amount of research in the field of face morphing. A very complete review of this field has been published by Scherhag *et al.* [156] in 2019 including both morphing techniques and also morphing attack detectors. Despite the large amount of publications, the research in this field is still in its infancy, with many open issues and challenges such as generating high-quality morphed images, the lack of metrics for reporting the vulnerability of face recognition systems to morphing attacks, etc. It is important to highlight the lack of publicly available databases and benchmarks what makes it difficult to perform a fair comparison among studies. In order to overcome this aspect, Raja *et al.* has recently presented an interesting framework for morphing attack detection [159], including a publicly available database, evaluation platform, and benchmark²⁸. The database comprises morphed and real images collected in three different sites constituting 1800 photographs of 150 subjects. Morphing images were generated using 6 different algorithms, presenting a wide variety of possible approaches.

Regarding the face morphing detectors, different approaches have been proposed in the literature based on different features, e.g.: the reduction of face details due to blending operations [160], Fourier spectrum of sensor pattern noise [161], differences between the facial landmarks [162,163], and pure deep learning features [164,165]. In addition, approaches based on face de-morphing have been studied in order to restore the accomplice's facial image [166,167].

7.2. Face de-Identification

The main goal of face de-identification (de-ID) is to remove the identity information present on a face image or video in order to preserve the privacy of the person [168]. This can be achieved in several ways.

²⁸ <https://biolab.csr.unibo.it/fvcongoing>.

The simplest way can be just to obfuscate the face by blurring or pixelation (e.g., in Google Maps Street View). More sophisticated methods try to provide face images with different identities but maintaining all other factors (pose, expression, illumination, etc.) unaltered. Therefore, the concept of face de-ID is very general. One possible option to achieve face de-identification could be through face identity swap.

Earlier works in this area were based on applying face de-ID to still images. In [169] Gross *et al.* presented a multi-factor framework for de-ID, which combined linear, bilinear, and quadratic models. They showed their method was able to protect privacy while preserving data utility on an expression-variant face database. Recently, the developments of image synthesis methods based on generative deep neural networks, in particular GAN, have inspired new face de-ID methods such as [170–175], which use synthesised faces to replace the original ones. Also, in [176], the authors proposed the use of Semi-Adversarial Networks (SAN) to confound arbitrary face-based gender classifiers.

More recently, in [177] Gafni *et al.* presented in 2019 a method that provides face de-ID with convincing performance even in unconstrained videos. Their approach is based on an adversarial autoencoder coupled with a trained face classifier. This way they can achieve a rich latent space, embedding both identity and expression information. Also, in [178] a new face de-ID method based on a deep transfer model was presented. This method treats the non-identity related facial attributes as the style of the original faces, and uses a trained facial attribute transfer model to extract and map them to different faces achieving very promising results both in single images and videos.

Some other related studies in this area work directly over face representations or deep face models by eliminating there undesired or protected information like identity, gender, or facial expressions [179–181]. Once that protected information has been disentangled, a face image or video can then be generated based on the new representations originated in which the protected information has been eliminated, reduced, or obfuscated.

7.3. Audio-to-Video and text-to-Video

A related topic to facial expression swap is the synthesis of video from audio or text. These types of video face manipulations are also known as lip-sinc deep fakes [182]. Popular examples can be seen on the Internet^{29,30}.

Regarding the synthesis of fake videos from audio (audio-to-video), Suwajanakorn *et al.* presented in [125] an approach to synthesise high quality videos of a person (Obama in this case) speaking with accurate lip sync. For this, they used as input to their approach many hours of previous videos of the person together with a new audio recording. In their approach they employed a recurrent neural network (based on LSTMs) to learn the mapping from raw audio features to mouth shapes. Then, based on the mouth shape at each frame, they synthesised high quality mouth texture, and composited it with 3D pose matching to create the new video to match the input audio track, producing photorealistic results.

In [183], Song *et al.* proposed an approach based on a novel conditional recurrent generation network that incorporates both image and audio features in the recurrent unit for temporal dependency, and also a pair of spatial-temporal discriminators for better image/video quality. As a result, their approach can model both lip and mouth together with expression and head pose variations as a whole, achieving much more realistic results. The source code is publicly available in GitHub³¹. Also, in [184] Song *et al.* presented a dynamic method not assuming a person-specific rendering network like in [125]. In their approach they are able to generate very realistic fake videos by carrying out a 3D face model reconstruction from the input video plus a recurrent network to translate

the source audio into expression parameters. Finally, they introduced a novel video rendering network and a dynamic programming method to construct a temporally coherent and photo-realistic video. Video results are shown on the Internet³².

Another interesting approach was presented in [185]. Zhou *et al.* proposed a novel framework called Disentangled Audio-Visual System (DAVS), which generates high quality talking face videos using disentangled audio-visual representation. Both audio and video speech information can be employed as input guidance. The source code is publicly available in GitHub³³.

Regarding the synthesis of fake videos from text (text-to-video), Fried *et al.* proposed in [186] a method that takes as input a video of a person speaking and the desired text to be spoken, and synthesises a new video in which the person's mouth is synchronised with the new words. In particular, their method automatically annotates an input talking-head video with phonemes, visemes, 3D face pose and geometry, reflectance, expression and scene illumination per frame. Finally, a recurrent video generation network creates a photorealistic video that matches the edited transcript. Examples of the fake videos generated with this approach are publicly available³⁴.

To the best of our knowledge, there are no publicly available databases and benchmarks related to audio- and text-to-video fake detection content. Research on this topic is usually carried out through the synthesis of in-house data using publicly available implementations like the ones described in this section.

Recent studies have analysed how easy is to detect audio- and text-to-video fake content. In [182], Agarwal *et al.* proposed a fake detection method that exploits the inconsistencies that exist between the dynamics of the mouth shape (visemes) and the spoken phoneme. They focused on some particular visemes in which the mouth must be completely closed and observed that this did not happen in many manipulated videos. Their proposed approach achieved good results, specially as the length of the video increases.

8. Concluding remarks

Motivated by the ongoing success of digital face manipulations, specially DeepFakes, this survey provides a comprehensive panorama of the field, including details of up-to-date: *i)* types of facial manipulations, *ii)* facial manipulation techniques, *iii)* public databases for research, and *iv)* benchmarks for the detection of each facial manipulation group, including key results achieved by the most representative manipulation detection approaches.

Generally speaking, most current face manipulations seem easy to be detected under controlled scenarios, i.e., when fake detectors are evaluated in the same conditions they are trained for. This fact has been demonstrated in most of the benchmarks included in this survey, achieving very low error rates in manipulation detection. However, this scenario may not be very realistic as fake images and videos are usually shared on social networks, suffering from high variations such as compression level, resizing, noise, etc. Also, facial manipulation techniques are continuously improving. These factors motivate further research on the generalisation ability of the fake detectors against unseen conditions. This aspect has been preliminary studied in different works [16,59–61]. Future research could be in the line of the latest publications [187,188] as they do not require fake videos for training, providing a better generalisation ability to unseen attacks.

Fusion techniques, at a feature or score level, could provide a better adaptation of the fake detectors to the different scenarios [189–191]. In fact, different fake detection approaches are already based on the combination of different sources of information, e.g., Zhou *et al.* proposed

²⁹ <https://www.youtube.com/watch?v=VWMEDacz3L4>.

³⁰ <https://www.bbc.com/news/technology-48607673>.

³¹ https://github.com/susanq/Talking_Face_Generation.

³² <https://wywu.github.io/projects/EBT/EBT.html>.

³³ <https://github.com/Hangz-nju-cuhk/Talking-Face-Generation-DAVS>.

³⁴ <https://www.ohadf.com/projects/text-based-editing/>.

in [95] a detection system based on the combination of steganalysis and pure deep learning features, whereas Rathgeb *et al.* proposed in [147] the combination of spatial and spectral features. Another two interesting fusion approaches have been recently presented in [192,193], combining RGB, Depth, and InfraRed information to detect physical face attacks. Also, face weighting approaches have been proposed in order to detect fake videos using multiple frames [194]. Finally, fusion of other sources of information such as the text, keystroke, or the audio that accompanies the videos when uploading them to social networks could be very valuable to improve the detectors [195–198].

In addition to the traditional fake detectors based only on the image/video information, novel schemes should be studied in order to provide more robust tools. One example of this is the work presented by Tursman *et al.* in [199]. The authors proposed to detect fake content via social verification at capture time: the arbiters of truthfulness are a group of video cameras that synchronously capture a speaker, collectively reach consensus, and then sign their videos in real time as “true”. Approaches like this could further protect media content from attacks.

We highlight next the key aspects to improve and future trends to follow for each facial manipulation group:

- **Face Synthesis:** current manipulations are usually based on GAN architectures such as StyleGAN, providing very realistic images. Nevertheless, most detectors can easily distinguish between real and fake images, achieving accuracies close to 100%. This is produced due to fake images are characterised by specific GAN fingerprints. But, what if we are able to remove those GAN fingerprints or add some noise patterns while keeping very realistic synthetic images? Recent approaches have focused on this research line, which represents a challenge even for the best manipulation detection systems [16,78,200].
- **Identity Swap:** although many different approaches have been proposed in the literature, it is certainly difficult to decide which is the best one. This is produced due to many different factors. First, most approaches are trained for a specific database and compression level, achieving in general very good results. However, they all show poor generalisation results to unseen conditions. In addition, the fact that different metrics (i.e., Acc., AUC, EER, etc.) and experimental protocols are usually considered does not help to achieve fair comparisons among studies. All these aspects should be further considered to advance in the field. Furthermore, we want to highlight the detection results achieved in the latest DeepFake databases of the 2nd generation such as DFDC and Celeb-DF [40,83]. While fake detectors already achieve AUC results close to 100% in databases of the 1st generation such as UADFV and FaceForensics++ [12,82], they all suffer from a high performance degradation on the latest ones, in particular for the Celeb-DF database with AUC results below 60% in most cases. Therefore, more efforts are needed to further improve current fake detection systems, for example, through large-scale challenges and benchmarks such as the recent DFDC³⁵.
- **Attribute Manipulation:** the same aspect highlighted for the face synthesis (GAN fingerprint removal) also applies here as most manipulations are based on GAN architectures. In addition, it is also interesting to remark the scarcity of public databases for research (only the DFFD database is publicly available [17]), and the lack of standard experimental protocols to perform fair comparisons among studies.
- **Expression Swap:** contrary to the identity swap, which has rapidly evolved with the release of improved DeepFake databases, the only public database in expression swap is FaceForensics++, to the best of our knowledge. This database is characterised by visual artifacts that are easy to detect, achieving therefore AUC results close to 100% in several fake detection approaches. We encourage researchers to

generate and make public more realistic databases based on recent techniques [125,151,184].

All these aspects, together with the development of improved GAN approaches and the recent DeepFake Detection Challenge (DFDC) will foster the new generation of realistic fake images/videos [70] together with more advanced techniques for face manipulation detection.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

CRediT authorship contribution statement

Ruben Tolosana: Conceptualization, Investigation, Writing - original draft, Writing - review & editing, Visualization, Funding acquisition. **Ruben Vera-Rodriguez:** Conceptualization, Writing - review & editing, Visualization, Funding acquisition. **Julian Fierrez:** Conceptualization, Writing - review & editing, Visualization, Funding acquisition. **Aythami Morales:** Conceptualization, Writing - review & editing, Visualization, Funding acquisition. **Javier Ortega-Garcia:** Conceptualization, Writing - review & editing, Visualization, Funding acquisition.

Acknowledgments

This work has been supported by projects: PRIMA (H2020-MSCA-ITN-2019-860315), TRESPASS-ETN (H2020-MSCA-ITN-2019-860813), BIBECA (MINECO/FEDER RTI2018-101248-B-I00), Bio-Guard (Ayudas Fundación BBVA a Equipos de Investigación Científica 2017), and Accenture. Ruben Tolosana is supported by Consejería de Educación, Juventud y Deporte de la Comunidad de Madrid y Fondo Social Europeo.

References

- [1] P. Korshunov, S. Marcel, Deepfakes: a New Threat to Face Recognition? Assessment and Detection, arXiv:1812.08685 (2018).
- [2] D. Citron, How DeepFake Undermine Truth and Threaten Democracy, 2019, URL <https://www.ted.com>.
- [3] R. Cellan-Jones, Deepfake Videos Double in Nine Months, 2019, URL <https://www.bbc.com/news/technology-49961089>.
- [4] BBC Bitesize, Deepfakes: What Are They and Why Would I Make One?, 2019, URL <https://www.bbc.co.uk/bitesize/articles/zfkwcqt>.
- [5] A. Swaminathan, M. Wu, K.J.R. Liu, Digital image forensics via intrinsic fingerprints, IEEE Trans. Inf. Forensics Secur. 3 (1) (2008) 101–117.
- [6] H. Farid, Image forgery detection, IEEE Signal Process. Mag. 26 (2) (2009) 16–25.
- [7] M. Stamm, K. Liu, Forensic detection of image manipulation using statistical intrinsic fingerprints, IEEE Trans. Inf. Forensics Secur. 5 (3) (2010) 492–506.
- [8] A. Rocha, W. Scheirer, T. Boult, S. Goldenstein, Vision of the unseen: current trends and challenges in digital image and video forensics, ACM Comput. Surv. 43 (4) (2011) 1–42.
- [9] S. Milani, M. Fontani, P. Bestagini, M. Barni, A. Piva, M. Tagliasacchi, S. Tubaro, An overview on video forensics, APSIPA Transactions on Signal and Information Processing 1 (2012) 1–18.
- [10] A. Piva, An overview on image forensics, ISRN Signal Processing 2013 (2013) 1–22.
- [11] P. Korus, Digital image integrity - a survey of protection and verification techniques, Digit. Signal Process. 71 (2017) 1–26.
- [12] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, M. Nießner, FaceForensics++: Learning to Detect Manipulated Facial Images, in: Proc. IEEE/CVF International Conference on Computer Vision, 2019.
- [13] J. Galbally, S. Marcel, J. Fierrez, Biometric anti-Spoofing methods: A Survey in face recognition, IEEE Access 2 (2014) 1530–1552.
- [14] A. Hadid, N. Evans, S. Marcel, J. Fierrez, Biometrics systems under spoofing attack: an evaluation methodology and lessons learned, IEEE Signal Process. Mag. 32 (5) (2015) 20–30.
- [15] S. Marcel, M. Nixon, J. Fierrez, N. Evans, Handbook of biometric anti-Spoofing (2nd edition), 2019.
- [16] J. Neves, R. Tolosana, R. Vera-Rodriguez, V. Lopes, H. Proença, J. Fierrez, GAN-printR: improved fakes and evaluation of the state-of-the-Art in face manipulation detection, IEEE J. Sel. Top. Signal Process. (2020).
- [17] H. Dang, F. Liu, J. Stehouwer, X. Liu, A. Jain, On the Detection of Digital Face Manipulation, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020.

³⁵ <https://deepfakedetectionchallenge.ai/>.

- [18] C. Canton, L. Davis, E. Delp, P. Flynn, S. McCloskey, L. Leal-Taixe, P. Natsev, C. Bregler, Applications of Computer Vision and Pattern Recognition to Media Forensics, in: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. URL <https://sites.google.com/view/mediaforensics2019>
- [19] B. Biggio, P. Korshunov, T. Mensink, G. Patrini, D. Rao, A. Sadhu, Synthetic Realities: Deep Learning for Detecting AudioVisual Fakes, in: *International Conference on Machine Learning*, 2019. URL <https://sites.google.com/view/audiovisualfakes-icml2019/>
- [20] L. Verdoliva, P. Bestagini, *Multimedia Forensics*, ACM Multimedia, 2019. URL <https://acmmm.org/tutorials/#tut3>
- [21] K. Raja, N. Damer, C. Chen, A. Dantcheva, A. Czajka, H. Han, R. Ramachandra, Workshop on Deepfakes and Presentation Attacks in Biometrics, in: *IEEE Winter Conference on Applications of Computer Vision*, 2020. URL <https://sites.google.com/view/wacv2020-deeppab/>
- [22] M. Barni, S. Battiato, G. Boato, H. Farid, N. Memon, *MultiMedia Forensics in the Wild*, in: *IEEE International Conference on Pattern Recognition*, 2020. URL <https://iplab.dmi.unict.it/mmforwild/>
- [23] C. Bregler, M. Covell, M. Slaney, Video rewrite: driving visual speech with audio, *Comput. Graph. (ACM)* 31 (2) (1997) 353–361.
- [24] D.P. Kingma, M. Welling, Auto-Encoding Variational Bayes, in: *Proc. International Conference on Learning Representations*, 2013.
- [25] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative Adversarial Nets, in: *Proc. Advances in Neural Information Processing Systems*, 2014.
- [26] I. Yerushalmy, H. Hel-Or, Digital image forgery detection based on lens and sensor aberration, *Int. J. Comput. Vis.* 92 (1) (2011) 71–91.
- [27] A.C. Popescu, H. Farid, Exposing digital forgeries in color filter array interpolated images, *IEEE Trans. Signal Process.* 53 (10) (2005) 3948–3959.
- [28] H. Cao, A.C. Kot, Accurate detection of demosaicing regularity for digital image forensics, *IEEE Trans. Inf. Forensics Secur.* 4 (4) (2009) 899–910.
- [29] Z. Lin, J. He, X. Tang, C. Tang, Fast, automatic and fine-grained tampered JPEG image detection via DCT coefficient analysis, *Pattern Recognit* 42 (11) (2009) 2492–2501.
- [30] Y.L. Chen, C.T. Hsu, Detecting recompression of JPEG images via periodicity analysis of compression artifacts for tampering detection, *IEEE Trans. Inf. Forensics Secur.* 6 (2) (2011) 396–406.
- [31] I. Amerini, L. Ballan, R. Caldelli, A. Bimbo, G. Serra, A SIFT-Based forensic method for copy-move attack detection and transformation recovery, *IEEE Trans. Inf. Forensics Secur.* 6 (3) (2011) 1099–1110.
- [32] D. Cozzolino, G. Poggi, L. Verdoliva, Splicebuster: A new blind image splicing detector, in: *Proc. IEEE International Workshop on Information Forensics and Security*, 2015, pp. 1–6.
- [33] A. Gironi, M. Fontani, T. Bianchi, A. Piva, M. Barni, A Video Forensic Technique for Detecting Frame Deletion and Insertion, in: *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, 2014, pp. 6226–6230.
- [34] Y. Wu, X. Jiang, T. Sun, W. Wang, Exposing Video Inter-Frame Forgery based on Velocity Field Consistency, in: *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, 2014, pp. 2674–2678.
- [35] B.C. Hosler, M.C. Stamm, Detecting Video Speed Manipulation, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [36] H. Allcott, M. Gentzkow, Social media and fake news in the 2016 election, *Journal of Economic Perspectives* 31 (2) (2017) 211–236.
- [37] D.M. Lazer, M.A. Baum, Y. Benkler, A.J. Berinsky, K.M. Greenhill, F. Menczer, M.J. Metzger, B. Nyhan, G. Pennycook, D. Rothschild, et al., The science of fake news, *Science* 359 (6380) (2018) 1094–1096.
- [38] I.M. J. Kietzmann L.W. Lee, T. Kietzmann, Deepfakes: trick or treat? *Bus. Horiz.* 63 (2) (2020) 135–146.
- [39] L. Verdoliva, Media forensics and DeepFakes: an overview, *IEEE J. Sel. Top. Signal Process.* (2020).
- [40] Y. Li, X. Yang, P. Sun, H. Qi, S. Lyu, Celeb-DF: A Large-Scale Challenging Dataset for DeepFake Forensics, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [41] T. Karras, S. Laine, T. Aila, A Style-Based Generator Architecture for Generative Adversarial Networks, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [42] E. Gonzalez-Sosa, J. Fierrez, R. Vera-Rodriguez, F. Alonso-Fernandez, Facial soft biometrics for recognition in the wild: recent works, annotation and COTS evaluation, *IEEE Trans. Inf. Forensics Secur.* 13 (8) (2018) 2001–2014.
- [43] Y. Choi, M. Choi, M. Kim, J. Ha, S. Kim, J. Choo, StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- [44] M. Liu, Y. Ding, M. Xia, X. Liu, E. Ding, W. Zuo, S. Wen, STGAN: A Unified Selective Transfer Network for Arbitrary Image Attribute Editing, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [45] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, M. Nießner, Face2face: Real-Time Face Capture and Reenactment of RGB Videos, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016.
- [46] J. Thies, M. Zollhofer, M. Nießner, Deferred neural rendering: image synthesis using neural textures, *ACM Trans. Graph.* 38 (66) (2019) 1–12.
- [47] 100,000 Faces Generated by AI, 2018, URL <https://generated.photos/>.
- [48] T. Karras, T. Aila, S. Laine, J. Lehtinen, Progressive Growing of GANs for Improved Quality, Stability, and Variation, in: *Proc. International Conference on Learning Representations*, 2018.
- [49] F. Marra, D. Gragnaniello, L. Verdoliva, G. Poggi, Do GANs Leave Artificial Fingerprints? in: *Proc. IEEE Conference on Multimedia Information Processing and Retrieval*, 2019, pp. 506–511.
- [50] M. Albright, S. McCloskey, Source Generator Attribution via Inversion? in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [51] A. Jain, P. Majumdar, R. Singh, M. Vatsa, Detecting GANs and Retouching based Digital Alterations via DAD-HCNN, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [52] Z. Liu, P. Luo, X. Wang, X. Tang, Deep Learning Face Attributes in the Wild, in: *Proc. IEEE/CVF International Conference on Computer Vision*, 2015.
- [53] D. Yi, Z. Lei, S. Liao, S. Li, Learning Face Representation From Scratch, *arXiv:1411.7923* (2014).
- [54] Q. Cao, L. Shen, W. Xie, O. Parkhi, A. Zisserman, VGGFace2: A Dataset for Recognising Faces Across Pose and Age, in: *Proc. International Conference on Automatic Face & Gesture Recognition*, 2018.
- [55] S. McCloskey, M. Albright, Detecting GAN-Generated Imagery Using Color Cues, *arXiv:1812.08247* (2018).
- [56] R. Wang, L. Ma, F. Juefei-Xu, X. Xie, J. Wang, Y. Liu, FakeSpotter: A Simple Baseline for Spotting AI-Synthesized Fake Faces, *arXiv:1909.06122* (2019).
- [57] L. Guarnera, O. Giudice, S. Battiato, DeepFake Detection by Analyzing Convolutional Traces, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [58] L. Nataraj, T. Mohammed, B. Manjunath, S. Chandrasekaran, A. Flenner, J. Bappy, A. Roy-Chowdhury, Detecting GAN generated fake images using co-Occurrence matrices, *Electronic Imaging* (5) (2019) 1–7.
- [59] N. Yu, L. Davis, M. Fritz, Attributing Fake Images to GANs: Analyzing Fingerprints in Generated Images, in: *Proc. IEEE/CVF International Conference on Computer Vision*, 2019.
- [60] F. Marra, C. Saltori, G. Boato, L. Verdoliva, Incremental Learning for the Detection and Classification of GAN-Generated Images, in: *Proc. IEEE International Workshop on Information Forensics and Security*, 2019.
- [61] N. Hulzebosch, S. Ibrahimi, M. Worringer, Detecting CNN-Generated Facial Images in Real-World Scenarios, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [62] H. Guan, M. Kozak, E. Robertson, Y. Lee, A. Yates, A. Delgado, D. Zhou, T. Kheyrkhan, J. Smith, J. Fiscus, MFC Datasets: Large-Scale Benchmark Datasets for Media Forensic Challenge Evaluation, in: *Proc. IEEE Winter Applications of Computer Vision Workshops*, 2019.
- [63] O. Parkhi, A. Vedaldi, A. Zisserman, Deep Face Recognition, in: *Proc. British Machine Vision Conference*, 2015.
- [64] B. Amos, B. Ludwiczuk, M. Satyanarayanan, OpenFace: A General-Purpose Face Recognition Library with Mobile Applications, *CMU School of Computer Science*, 2016.
- [65] F. Schroff, D. Kalenichenko, J. Philbin, FaceNet: A Unified Embedding for Face Recognition and Clustering, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2015.
- [66] Y. Shen, J. Gu, X. Tang, B. Zhou, Interpreting the Latent Space of GANs for Semantic Face Editing, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [67] T.K. Moon, The expectation-Maximization algorithm, *IEEE Signal Process. Mag.* 13 (6) (1996) 47–60.
- [68] Z. He, W. Zuo, M. Kan, S. Shan, X. Chen, Attgan: facial attribute editing by only changing what you want, *IEEE Trans. Image Process.* 28 (11) (2019) 5464–5478.
- [69] W. Cho, S. Choi, D.K. Park, I. Shin, J. Choo, Image-to-Image Translation via Group-Wise Deep Whitening and Coloring Transformation, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [70] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, T. Aila, Analyzing and Improving the Image Quality of StyleGAN, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [71] J. Zhu, T. Park, P. Isola, A. Efros, Unpaired Image-to-Image Translation Using Cycle-Consistent Adversarial Networks, in: *Proc. IEEE/CVF International Conference on Computer Vision*, 2017.
- [72] T. Miyato, T. Kataoka, M. Koyama, Y. Yoshida, Spectral Normalization for Generative Adversarial Networks, in: *Proc. International Conference on Learning Representations*, 2018.
- [73] M. Bellemare, I. Danihelka, W. Dabney, S. Mohamed, B. Lakshminarayanan, S. Hoyer, R. Munos, The Cramer Distance as a Solution to Biased Wasserstein Gradients, *arXiv:1705.10743* (2017).
- [74] M. Binkowski, D. Sutherland, M. Arbel, A. Gretton, Demystifying MMD GANs, in: *Proc. International Conference on Learning Representations*, 2018.
- [75] S. Rebuffi, A. Kolesnikov, G. Sperl, C. Lampert, iCaRL: Incremental Classifier and Representation Learning, in: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017.
- [76] D. Kingma, P. Dhariwal, Glow: Generative Flow with Invertible 1x1 Convolutions, in: *Proc. Advances in Neural Information Processing Systems*, 2018.
- [77] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, I. Polosukhin, Attention is All You Need, in: *Proc. Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [78] D. Cozzolino, J. Thies, A. Rössler, M. Nießner, L. Verdoliva, SpoC: Spoofing Camera Fingerprints, *arXiv:1911.12069* (2019).
- [79] D. Cozzolino, J. Thies, A. Rössler, C. Riess, M. Nießner, L. Verdoliva, ForensicTransfer: Weakly-Supervised Domain Adaptation for Forgery Detection, *arXiv:1812.02510* (2018).

- [80] M. Huh, A. Liu, A. Owens, A. Efros, Fighting Fake News: Image Splice Detection Via Learned Self-Consistency, in: Proc. European Conference on Computer Vision, 2018.
- [81] P. Zhou, X. Han, V. Morariu, L. Davis, Learning Rich Features for Image Manipulation Detection, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018.
- [82] Y. Li, M. Chang, S. Lyu, In Ictu Oculi: Exposing AI Generated Fake Face Videos by Detecting Eye Blinking, in: Proc. IEEE International Workshop on Information Forensics and Security, 2018.
- [83] B. Dolhansky, R. Howes, B. Pfaff, N. Baram, C.C. Ferrer, The Deepfake Detection Challenge (DFDC) Preview Dataset, arXiv:1910.08854 (2019).
- [84] Google AI, Contributing Data to Deepfake Detection Research, 2019, URL <https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html>.
- [85] C. Sanderson, B. Lovell, Multi-Region Probabilistic Histograms for Robust and Scalable Identity Inference, in: Proc. International Conference on Biometrics, 2009.
- [86] K. Zhang, Z. Zhang, Z. Li, Y. Qiao, Joint face detection and alignment using multitask cascaded convolutional networks, IEEE Signal Process. Lett. 23 (10) (2016) 1499–1503.
- [87] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, M. Nießner, FaceForensics: A Large-Scale Video Dataset for Forgery Detection in Human Faces, arXiv:1803.09179 (2018).
- [88] P. Pérez, M. Gangnet, A. Blake, Poisson image editing, ACM Trans. Graph. 22 (3) (2003) 313–318.
- [89] F. Matern, C. Riess, M. Stamminger, Exploiting Visual Artifacts to Expose DeepFakes and Face Manipulations, in: Proc. IEEE Winter Applications of Computer Vision Workshops, 2019.
- [90] X. Yang, Y. Li, S. Lyu, Exposing Deep Fakes Using Inconsistent Head Poses, in: Proc. International Conference on Acoustics, Speech and Signal Processing, 2019.
- [91] S. Agarwal, H. Farid, Protecting World Leaders Against Deep Fakes, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2019.
- [92] T. Jung, S. Kim, K. Kim, Deepvision: deepfakes detection using human eye blinking pattern, IEEE Access 8 (2020) 83144–83154.
- [93] Y. Li, S. Lyu, Exposing DeepFake Videos By Detecting Face Warping Artifacts, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2019.
- [94] D. Afchar, V. Nozick, J. Yamagishi, I. Echizen, MesoNet: a Compact Facial Video Forgery Detection Network, in: Proc. IEEE International Workshop on Information Forensics and Security, 2018.
- [95] P. Zhou, X. Han, V. Morariu, L. Davis, Two-Stream Neural Networks for Tampered Face Detection, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2017.
- [96] H. Nguyen, F. Fang, J. Yamagishi, I. Echizen, Multi-task Learning For Detecting and Segmenting Manipulated Facial Images and Videos, arXiv:1906.06876 (2019).
- [97] H.H. Nguyen, J. Yamagishi, I. Echizen, Use of a Capsule Network to Detect Fake Images and Videos, arXiv:1910.12467 (2019).
- [98] E. Sabir, J. Cheng, A. Jaiswal, W. AbdAlmageed, I. Masi, P. Natarajan, Recurrent Convolutional Strategies for Face Manipulation Detection in Videos, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2019.
- [99] D. Güera, E. Delp, Deepfake Video Detection Using Recurrent Neural Networks, in: Proc. International Conference on Advanced Video and Signal Based Surveillance, 2018.
- [100] R. Tolosana, S. Romero-Tapiador, J. Fierrez, R. Vera-Rodriguez, DeepFakes Evolution: Analysis of Facial Regions and Fake Detection Performance, arXiv:2004.07532 (2020).
- [101] P. Korshunov, S. Marcel, Speaker Inconsistency Detection in Tampered Video, in: Proc. European Signal Processing Conference, 2018.
- [102] J. Galbally, S. Marcel, J. Fierrez, Image quality assessment for fake biometric detection: application to iris, fingerprint and face recognition, IEEE Trans. Image Process. 23 (2) (2014) 710–724.
- [103] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, 2016.
- [104] D. King, DLib-ML: A Machine Learning toolkit, Journal of Machine Learning Research 10 (2009) 1755–1758.
- [105] T. Baltrušaitis, A. Zadeh, Y. Lim, L. Morency, OpenFace 2.0: Facial Behavior Analysis Toolkit, in: Proc. International Conference on Automatic Face & Gesture Recognition, 2018.
- [106] R. Daza, A. Morales, J. Fierrez, R. Tolosana, mEBAL: A Multimodal Database for Eye Blink Detection and Attention Level Estimation, arXiv:2006.05327 (2020).
- [107] R. Ranjan, V.M. Patel, R. Chellappa, Hyperface: a deep multi-Task learning framework for face detection, landmark localization, pose estimation, and gender recognition, IEEE Trans. Pattern Anal. Mach. Intell. 41 (1) (2017) 121–135.
- [108] T. Soukupová, J. Cech, Eye Blink Detection Using Facial Landmarks, in: Proc. Computer Vision Winter Workshop, 2016.
- [109] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, arXiv:1409.1556 (2014).
- [110] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.
- [111] K. He, X. Zhang, S. Ren, J. Sun, Spatial pyramid pooling in deep convolutional networks for visual recognition, IEEE Trans. Pattern Anal. Mach. Intell. 37 (9) (2015) 1904–1916.
- [112] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going Deeper with Convolutions, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2015.
- [113] D. Cozzolino, G. Poggi, L. Verdoliva, Recasting Residual-Based Local Descriptors as Convolutional Neural Networks: an Application to Image Forgery Detection, in: Proc. ACM Workshop on Information Hiding and Multimedia Security, 2017.
- [114] B. Bayar, M. Stamm, A Deep Learning Approach to Universal Image Manipulation Detection Using a New Convolutional Layer, in: Proc. ACM Workshop on Information Hiding and Multimedia Security, 2016.
- [115] N. Rahmouni, V. Nozick, J. Yamagishi, I. Echizen, Distinguishing Computer Graphs from Natural Images Using Convolution Neural Networks, in: Proc. IEEE Workshop on Information Forensics and Security, 2017.
- [116] F. Chollet, Xception: Deep Learning with Depthwise Separable Convolutions, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2017.
- [117] J. Deng, W. Dong, R. Socher, L. Li, K. Li, L. Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2009.
- [118] G.E. Hinton, A. Krizhevsky, S. Krizhevsky, S.D. Wang, Transforming Auto-Encoders, in: International Conference on Artificial Neural Networks, 2011, pp. 44–51.
- [119] S. Sabour, N. Frosst, G.E. Hinton, Dynamic Routing Between Capsules, in: Proc. Advances in Neural Information Processing Systems, 2017, pp. 3856–3866.
- [120] G.E. Hinton, S. Sabour, N. Frosst, Matrix Capsules with EM routing, in: Proc. International Conference on Learning Representations Workshop, 2018.
- [121] Y. Wang, A. Dantcheva, A Video is Worth More than 1000 Lies. Comparing 3DCNN Approaches for Detecting Deepfakes, in: Proc. IEEE International Conference on Automatic Face and Gesture Recognition, 2020.
- [122] J. Carreira, A. Zisserman, Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2017.
- [123] K. Hara, H. Kataoka, Y. Satoh, Can Spatiotemporal 3D CNNs Retrace the History of 2D CNNs and Imagenet? in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018.
- [124] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the Inception Architecture for Computer Vision, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2016.
- [125] S. Suwajanakorn, S. Seitz, I. Kemelmacher-Shlizerman, Synthesizing obama: learning lip sync from audio, ACM Trans. Graph. 36 (4) (2017) 1–13.
- [126] P. Isola, J. Zhu, T. Zhou, A. Efros, Image-to-Image Translation with Conditional Adversarial Networks, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2017.
- [127] J. Zhu, R. Zhang, D. Pathak, T. Darrell, A. Efros, O. Wang, E. Shechtman, Toward Multimodal Image-to-Image Translation, in: Proc. Advances in Neural Information Processing Systems, 2017.
- [128] T. Kim, M. Cha, H. Kim, J. Lee, J. Kim, Learning to Discover Cross-Domain Relations with Generative Adversarial Networks, in: Proc. International Conference on Machine Learning, 2017.
- [129] D. Bau, J. Zhu, H. Strobelt, B. Zhou, J. Tenenbaum, W. Freeman, A. Torralba, GAN Dissection: Visualizing and Understanding Generative Adversarial Networks, arXiv:1811.10597 (2018).
- [130] G. Perarnau, J.V.D. Weijer, B. Raducanu, J. Álvarez, Invertible Conditional GANs for Image Editing, in: Proc. Advances in Neural Information Processing Systems Workshops, 2016.
- [131] M. Li, W. Zuo, D. Zhang, Deep Identity-Aware Transfer of Facial Attributes, arXiv:1610.05586 (2016).
- [132] W. Shen, R. Liu, Learning Residual Images for Face Attribute Manipulation, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2017.
- [133] G. Lample, N. Zeghidour, N. Usunier, A. Bordes, L. Denoyer, M. Ranzato, Fader Networks: Manipulating Images by Sliding Attributes, in: Proc. Advances in Neural Information Processing Systems, 2017.
- [134] T. Xiao, J. Hong, J. Ma, ELEGANT: Exchanging Latent Encodings with GAN for Transferring Multiple Face Attributes, in: Proc. European Conference on Computer Vision, 2018.
- [135] M. Mirza, S. S. Osindero, Conditional Generative Adversarial Nets, arXiv:1411.1784 (2014).
- [136] J. Kim, J. Choi, J. Yi, M. Turk, Effective representation using ICA for face recognition robust to local distortion and partial occlusion, IEEE Trans. Pattern Anal. Mach. Intell. 12 (2005) 1977–1981.
- [137] A. Dantcheva, C. Chen, A. Ross, Can Facial Cosmetics Affect the Matching Accuracy of Face Recognition Systems? in: Proc. International Conference on Biometrics: Theory, Applications and Systems, 2012, pp. 391–398.
- [138] N. Kose, L. Aprville, J. Dugelay, Facial Makeup Detection Technique based on Texture and Shape Analysis, in: Proc. International Conference and Workshops on Automatic Face and Gesture Recognition, 2015.
- [139] P. Majumdar, A. Agarwal, R. Singh, M. Vatsa, Evading Face Recognition via Partial Tampering of Faces, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2019.
- [140] C. Rathgeb, A. Dantcheva, C. Busch, Impact and detection of facial beautification in face recognition: an overview, IEEE Access 7 (2019) 152667–152678.
- [141] C. Rathgeb, C.I. Satnoianu, N.E. Haryanto, K. Bernardo, C. Busch, Differential detection of facial retouching: a multi-Biometric approach, IEEE Access 8 (2020) 106373–106385.
- [142] A. Bharati, R. Singh, M. Vatsa, K. Bowyer, Detecting facial retouching using supervised deep learning, IEEE Trans. Inf. Forensics Secur. 11 (9) (2016) 1903–1913.
- [143] A. Jain, R. Singh, M. Vatsa, On Detecting GANs and Retouching based Synthetic Alterations, in: Proc. International Conference on Biometrics Theory, Applications and Systems, 2018.
- [144] S. Tariq, S. Lee, H. Kim, Y. Shin, S. Woo, Detecting Both Machine and Human Created Fake Face Images in the Wild, in: Proc. International Workshop on Multimedia Privacy and Security, 2018, pp. 81–87.

- [145] S. Wang, O. Wang, A. Owens, R. Zhang, A. Efros, Detecting Photoshopped Faces by Scripting Photoshop, in: Proc. IEEE/CVF International Conference on Computer Vision, 2019.
- [146] X. Zhang, S. Karaman, S. Chang, Detecting and Simulating Artifacts in GAN Fake Images, in: Proc. IEEE International Workshop on Information Forensics and Security, 2019.
- [147] C. Rathgeb, A. Botaljov, F. Stockhardt, S. Isadskiy, L. Debiase, A. Uhl, C. Busch, PRNU-Based detection of facial retouching, IET Biom. (2020).
- [148] P. Flynn, K. Bowyer, P. Phillips, Assessment of Time Dependency in Face Recognition: An Initial Study, in: Proc. International Conference on Audio-and Video-Based Biometric Person Authentication, 2003.
- [149] Adjust and Exaggerate Facial Features. Adobe Photoshop, 2016, URL <https://helpx.adobe.com/photoshop/how-to/face-aware-liquify.html>.
- [150] D. Sun, X. Yang, M.Y. Liu, J. Kautz, PWC-Net: CNNs for Optical Flow Using Pyramid, Warping, and Cost Volume, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018.
- [151] H. Averbuch-Elor, D. Cohen-Or, J. Kopf, M.F. Cohen, Bringing portraits to life, ACM Trans Graph 36 (6) (2017) 196.
- [152] E. Zakharov, A. Shysheya, E. Burkov, V. Lempitsky, Few-Shot Adversarial Learning of Realistic Neural Talking Head Models, in: Proc. IEEE/CVF International Conference on Computer Vision, 2019.
- [153] D. Zhu, S. Liu, W. Jiang, C. Gao, T. Wu, G. Guo, UGAN: Untraceable GAN for Multi-Domain Face Translation, arXiv:1907.11418 (2019).
- [154] I. Amerini, L. Galteri, R. Caldelli, A. Bimbo, Deepfake Video Detection through Optical Flow based CNN, in: Proc. IEEE/CVF International Conference on Computer Vision, 2019.
- [155] G. Wolberg, Image morphing: a survey, Vis. Comput. 14 (8–9) (1998) 360–372.
- [156] U. Scherhag, C. Rathgeb, J. Merkle, R. Breithaupt, C. Busch, Face recognition systems under morphing attacks: asurvey, IEEE Access 7 (2019) 23012–23026.
- [157] M. Gomez-Barrero, C. Rathgeb, U. Scherhag, C. Busch, Is Your Biometric System Robust to Morphing Attacks? in: Proc. IEEE International Workshop on Biometrics and Forensics, 2017.
- [158] P. Korshunov, S. Marcel, Vulnerability of Face Recognition to Deep Morphing, arXiv:1910.01933 (2019).
- [159] K. Raja, M. Ferrara, A. Franco, L. Spreeuwiers, I. Batskos, F. de Wit, M. Gomez-Barrero, U. Scherhag, D. Fischer, S. Venkatesh, J.M. Singh, G. Li, L. Bergeron, S. Isadskiy, R. Ramachandra, C. Rathgeb, D. Frings, U. Seidel, F. Knopjes, R. Veldhuis, D. Maltoni, C. Busch, Morphing Attack Detection - Database, Evaluation Platform and Benchmarking, arXiv:2006.06458 (2020).
- [160] C. Kraetzer, A. Makrushin, T. Neubert, M. Hildebrandt, J. Dittmann, Modeling Attacks on Photo-ID Documents and Applying Media Forensics for the Detection of Facial Morphing, in: Proc. ACM Workshop on Information Hiding and Multimedia Security, 2017, pp. 21–32.
- [161] L.B. Zhang, F. Peng, M. Long, Face Morphing Detection Using Fourier Spectrum of Sensor Pattern Noise, in: Proc. International Conference on Multimedia and Expo, 2018.
- [162] U. Scherhag, D. Budhrani, M. Gomez-Barrero, C. Busch, Detecting Morphed Face Images Using Facial Landmarks, in: Proc. IEEE International Conference on Image and Signal Processing, 2018.
- [163] N. Damer, V. Bolle, Y. Wainakh, F. Boutros, P. Terhöst, A. Braun, A. Kuijper, Detecting Face Morphing Attacks by Analyzing the Directed Distances of Facial Landmarks Shifts, in: German Conference on Pattern Recognition, 2018.
- [164] M. Ferrara, A. Franco, D. Maltoni, Face Morphing Detection in the Presence of Printing/Scanning and Heterogeneous Image Sources, arXiv:1901.08811 (2019).
- [165] U. Scherhag, C. Rathgeb, J. Merkle, C. Busch, Deep Face Representations for Differential Morphing Attack Detection, arXiv:2001.01202 (2020).
- [166] M. Ferrara, A. Franco, D. Maltoni, Face demorphing, IEEE Trans. Inf. Forensics Secur. 13 (4) (2017) 1008–1017.
- [167] F. Peng, L.B. Zhang, M. Long Min, FD-GAN: Face de-Morphing generative adversarial network for restoring Accomplice's facial image, IEEE Access 7 (2019) 75122–75131.
- [168] R. Gross, L. Sweeney, F. De la Torre, S. Baker, Model-Based Face De-Identification, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop, 2006.
- [169] R. Gross, L. Sweeney, J. Cohn, F. De la Torre, S. Baker, Face De-Identification, in: Protecting Privacy in Video Surveillance, Springer, 2009, pp. 129–146.
- [170] B. Meden, R.C. Malli, S. Fabijan, H.K. Ekenel, V. Štruc, P. Peer, Face deidentification with generative deep neural networks, IET Signal Proc. 11 (9) (2017) 1046–1054.
- [171] K. Brkic, I. Sikiric, T. Hrkac, Z. Kalafatic, I Know That Person: Generative Full Body and Face De-identification of People in Images, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2017.
- [172] Q. Sun, L. Ma, S.O. Joon, L.V. Gool, B. Schiele, M. Fritz, Natural and Effective Obfuscation by Head Inpainting, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018.
- [173] B. Meden, V. Emeršič, V. Štruc, P. Peer, K-Same-Net: k-Anonymity with generative deep neural networks for face deidentification, Entropy 20 (1) (2018) 60.
- [174] S. Guo, S. Feng, Y. Li, S. An, H. Dong, Integrating Diversity into Neural-Network-Based Face Deidentification, in: Proc. Chinese Control Conference, 2018.
- [175] Y.L. Pan, M.J. Huang, K.T. Ding, J.L. Wu, J.S. Jang, K-Same-Siamese-GAN: K-Same Algorithm with Generative Adversarial Network for Facial Image De-identification with Hyperparameter Tuning and Mixed Precision Training, in: Proc. International Conference on Advanced Video and Signal Based Surveillance, 2019.
- [176] V. Mirjalili, S. Raschka, A. Ross, Flowsan: privacy-Enhancing semi-Adversarial networks to confound arbitrary face-based gender classifiers, IEEE Access 7 (2019) 99735–99745.
- [177] O. Gafni, L. Wolf, Y. Taigman, Live Face De-Identification in Video, arXiv:1911.08348 (2019).
- [178] Y. Li, S. Lyu, De-Identification Without Losing Faces, in: Proc. of the ACM Workshop on Information Hiding and Multimedia Security, 2019.
- [179] M. Alvi, A. Zisserman, C. Nellaker, Turning a Blind Eye: Explicit Removal of Biases and Variation from Deep Neural Network Embeddings, in: Proc. European Conference on Computer Vision, 2018.
- [180] A. Morales, J. Fierrez, R. Vera-Rodriguez, SensitiveNets: Learning Agnostic Representations with Application to Face Recognition, arXiv:1902.00334 (2019).
- [181] S. Gong, X. Liu, A. Jain, DebFace: De-biasing Face Recognition, arXiv:1911.08080 (2019).
- [182] S. Agarwal, H. Farid, O. Fried, M. Agrawala, Detecting Deep-Fake Videos from Phoneme-Viseme Mismatches, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [183] Y. Song, J. Zhu, D. Li, A. Wang, H. Qi, Talking Face Generation by Conditional Recurrent Adversarial Network, in: Proc. International Joint Conference on Artificial Intelligence, 2019.
- [184] L. Song, W. Wu, C. Qian, R. He, C. Loy, Everybody's Talkin': Let Me Talk as You Want, arXiv:2001.05201 (2020).
- [185] H. Zhou, Y. Liu, Z. Liu, P. Luo, X. Wang, Talking Face Generation by Adversarially Disentangled Audio-Visual Representation, in: Proc. AAAI Conference on Artificial Intelligence, 2019.
- [186] O. Fried, A. Tewari, M. Zollhöfer, A. Finkelstein, E. Shechtman, D.B. Goldman, K. Genova, Z. Jin, C. Theobalt, M. Agrawala, Text-Based editing of talking-Head video, ACM Trans. Graph. 38(4) 1–14.
- [187] H. Khalid, S.S. Woo, OC-FakeDect: Classifying Deepfakes Using One-class Variational Autoencoder, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [188] S. Fernandes, S. Raj, R. Ewert, J.S. Pannu, S.K. Jha, E. Ortiz, I. Vintila, M. Salter, Detecting Deepfake Videos using Attribution-Based Confidence Metric, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [189] J. Fierrez, A. Morales, R. Vera-Rodriguez, D. Camacho, Multiple classifiers in biometrics. part 1: fundamentals and review, Information Fusion 44 (2018) 57–64.
- [190] J. Fierrez, A. Morales, R. Vera-Rodriguez, D. Camacho, Multiple classifiers in biometrics. part 2: trends and challenges, Information Fusion 44 (2018) 103–112.
- [191] R.S. M. Singh, A. Ross, A comprehensive overview of biometric fusion, Information Fusion 52 (2019) 187–205.
- [192] Q. Yang, X. Zhu, J.K. Fwu, Y. Ye, G. You, Y. Zhu, PipeNet: Selective Modal Pipeline of Fusion Network for Multi-Modal Face Anti-Spoofing, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [193] Z. Yu, Y. Qin, X. Li, Z. Wang, C. Zhao, Z. Lei, G. Zhao, Multi-Modal Face Anti-Spoofing Based on Central Difference Networks, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [194] D.M. Montserrat, H. Hao, S.K. Yarlagadda, S. Baireddy, R. Shao, J. Horvath, E. Bartusiak, J. Yang, D. Güera, F. Zhu, E.J. Delp, Deepfakes Detection with Automatic Face Weighting, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [195] T. Agrawal, R. Gupta, S. Narayanan, Multimodal Detection of Fake Social Media Use through a Fusion of classification and Pairwise Ranking Systems, in: Proc. European Signal Processing Conference, 2017, pp. 1045–1049.
- [196] K. Shu, A. Sliva, S. Wang, J. Tang, H. Liu, Fake news detection on social media: A Data mining perspective, ACM SIGKDD Explorations Newsletter 19 (1) (2017) 22–36.
- [197] K. Shu, D. Mahudeswaran, H. Liu, Fakenewstracker: a tool for fake news collection, detection, and visualization, Comput Math Organ Theory 25 (1) (2019) 60–71.
- [198] A. Morales, A. Acien, J. Fierrez, J.V. Monaco, R. Tolosana, R. Vera-Rodriguez, J. Ortega-Garcia, Keystroke Biometrics in Response to Fake News Propagation in a Global Pandemic, in: Proc. IEEE Computer Software and Applications Conference Workshops, 2020.
- [199] E. Tursman, M. George, S. Kamara, J. Tompkin, Towards Untrusted Social Video Verification to Combat Deepfakes via Face Geometry Consistency, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.
- [200] N. Carlini, H. Farid, Evading Deepfake-Image Detectors with White- and Black-Box Attacks, in: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020.