
A SURVEY ON ANTI-SPOOFING METHODS FOR FACE RECOGNITION WITH RGB CAMERAS OF GENERIC CONSUMER DEVICES

Zuheng MING¹, Muriel VISANI^{1,2,*}, Muhammad Muzzamil LUQMAN¹, Jean-Christophe BURIE¹

¹ L3i Laboratory, La Rochelle University, France

² School of Information & Communication Technology, Hanoi University of Science and Technology, Vietnam
{zuheng.ming, muriel.visani, mluqma01, jcburie}@univ-lr.fr

ABSTRACT

The widespread deployment of face recognition-based biometric systems has made face Presentation Attack Detection (face anti-spoofing) an increasingly critical issue. This survey thoroughly investigates the face Presentation Attack Detection (PAD) methods, that only require RGB cameras of generic consumer devices, over the past two decades. We present an attack scenario-oriented typology of the existing face PAD methods and we provide a review of over 50 of the most recent face PAD methods and their related issues. We adopt a comprehensive presentation of the methods that have most influenced face PAD following the proposed typology, and in chronological order. By doing so, we depict the main challenges, evolutions and current trends in the field of face PAD, and provide insights on its future research. From an experimental point of view, this survey paper provides a summarized overview of the available public databases and extensive comparative experimental results of different PAD methods.

Keywords Biometrics, face recognition, face anti-spoofing, face Presentation Attack Detection (PAD), RGB camera-based anti-spoofing methods, deep learning, survey, computer vision, pattern recognition

1 Introduction

1.1 Background

In the past two decades, the advancement of technology in electronics and computer science has provided access to top-level technology devices at affordable prices to an important proportion of the world population. Various biometric systems have been widely deployed in real-life applications, such as on-line payment and e-commerce security, smartphone-based authentication, secured access control, biometric passport and border checks. Face recognition is among the most studied biometric technologies since the 90s [1], mainly for its numerous advantages compared to other biometrics. Indeed, faces are highly distinctive among individuals and face recognition can be implemented even in non-intrusive acquisition scenarios, or from a distance.

Recently, deep learning has dramatically improved the state-of-the-art performance of many computer vision tasks, such as image classification and object recognition [2, 3, 4]. With these significant progresses, face recognition has also made great breakthroughs such as the success of DeepFace [5], DeepIDs[6], VGG Face [7], FaceNet [8], SphereFace [9] and ArcFace [10]. One of these spectacular breakthroughs occurred between 2014 and 2015, when multiple groups [5, 11, 8] approached and then surpassed human-level recognition

*Correspondence: muriel.visani@univ-lr.fr.

accuracy on very challenging face benchmarks, such as LFW [12] or YTF [13]. Thanks to their convenience, excellent performances and great security levels, face recognition systems are among the most widespread biometric systems in the market, compared to other biometrics such as iris and fingerprints [14].

However, given face authentication systems' popularity, they became primary targets of Presentation Attacks (PAs) [15]. PAs are performed by malicious or ill-intentioned users who either aim at impersonating someone else's identity (impersonation attack), or at avoiding being recognized by the system (obfuscation attack). Yet, compared to face recognition performances, the vulnerabilities of face authentication systems to PAs have been much less studied.

The main objective of this paper is to present a detailed review of face PAD methods, that are crucial for assessing the vulnerability / robustness of current face recognition-based systems towards ill-intentioned users. Given the prevalence of biometric applications based on face authentication, such as online payment, it is crucial to protect genuine users against impersonation attacks in real-life scenarios. In this survey paper, we will focus more on impersonation detection. However, at the end of the paper, we will discuss obfuscation detection as well.

The next section provides a categorization of face PAs. Based on this categorization, we will present later in this paper a typology of existing face PAD methods and then a comprehensive review of such methods, with an extensive experimental comparison of these methods.

1.2 Categories of face Presentation Attacks

One can consider that there are basically two types of Presentation Attacks (PAs).

First, with the advent of internet and social medias where more and more people share photos or videos of their faces, such documents can be used by impostors to try and fool face authentication systems, for impersonation purposes. Such attacks are also called *impersonation (spoofing)* attacks.

Second, another (less studied) type of Presentation Attacks is called *obfuscation* attacks, where a person uses tricks to avoid being recognized by the system (but not necessarily by impersonating a legitimate user's identity).

In short, while impersonation (spoofing) attacks are generally performed by impostors who are willing to impersonate a legitimate user, obfuscation attacks aim at ensuring the user remains under the radar of the face recognition system. Despite their totally different objectives, both types of attacks are listed in the ISO standard [16] dedicated to biometric PAD.

In this survey paper, we focus on impersonation (spoofing) attacks, where the impostor might either use directly biometric data from a legitimate user to mount an attack, or to create Presentation Attack Instruments (PAIs, usually spoofs or fakes) that will be used for attacking the face recognition system.

Common PAs / PAIs can generally be categorized as photo attacks, video replay attacks, and 3D mask attacks (see Figure 1 for their categorization and Figure 2 for illustrations), whereas obfuscation attacks generally rely on tricks to hide the user's real identity, such as facial makeup, plastic surgery or face region occlusion.

Photo attacks (sometimes also called print attacks in the literature), and video replay attacks, are the most common attacks, due to the ever-increasing flow of face images available on the internet and prevalence of low-cost but high-resolution digital devices. Impostors can simply collect and re-use face samples of genuine users. Photo attacks are carried out by presenting to the face authentication system a picture of a genuine user. Several strategies are usually used by the impostors. Printed photo attacks (see Figure 2(a)) consist in presenting a picture printed on a paper (*e.g.*, A3/A4 paper, copper paper or professional photographic paper). On the other hand, in photo display attacks, the picture is displayed on the screen of a digital device such as a smartphone, a tablet or a laptop and then presented to the system. Moreover, as illustrated in Figure 2(b), printed photos can be warped (along a vertical and/or horizontal axis) to give some depth to the photo (this strategy is called warped photo attack). Cut photo attacks consist in using the picture as a photo mask where the mouth, eyes and/or nose regions have been cut out to introduce some liveness cues from the impostor's face behind the photo, such as eye blinking or mouth movement (see Figure 2(c)).

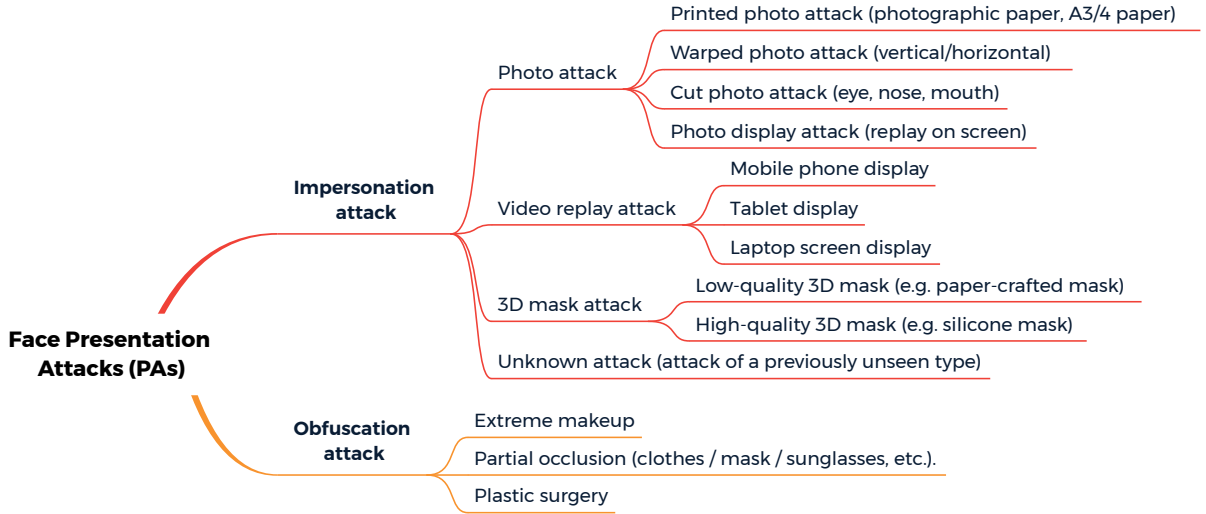


Figure 1: A typology of face Presentation Attacks (PAs).

Compared to static photo attacks, video replay attacks (see Figure 2(d)) are more sophisticated, as they introduce intrinsic dynamic information such as eye blinking, mouth movements, and changes in facial expressions, in order to mimic liveness [21].

Contrary to photo attacks or video replay attacks (that are generally 2D planar attacks, except for warped photo attacks), 3D mask attacks reconstruct 3D face artifacts. One can distinguish between low-quality 3D masks (*e.g.* crafted from a printed photo as illustrated in figure 2(e)) and high-quality 3D masks (*e.g.* made out of silicone, see Figure 2(f)). The high realism of the "face-like" 3D structure and the vivid imitation of human skin's texture of high-quality 3D masks makes it more challenging to detect 3D mask spoofing by traditional PAD methods (*i.e.* methods conceived to detect photo or video replay attacks [22, 23]). Nowadays, manufacturing a high-quality 3D mask is still expensive [24], complex, and relies on a complete 3D acquisition, generally requiring the user's cooperation [25]. Thus, 3D mask attacks are still far less frequent than photo or video replay attacks. However, with the popularization of 3D acquisition sensors, 3D mask attacks are expected to become more and more frequent in the coming years.

PAD methods for previously unseen attacks (unknown attacks) will be reviewed in Section 2.6: "New trends", as most of them are still under development and rely on recent approaches such as zero/few-shot learning.

Obfuscation attacks, whose objective is quite different from impersonation attacks (as the aim for the attacker is to remain unrecognized by the system), generally rely on facial makeup, plastic surgery or face region occlusion (*e.g.* using accessories such as scarves or sunglasses). However, in some cases, obfuscation attacks can also rely on the use of another person's biometric data. It fundamentally differs from usual spoofing attacks in its primary objective. However, in some cases, the PAIs for obfuscation attacks can be similar to the ones used for impersonation attacks: *e.g.* the face mask of another person. While most of the PAD methods reviewed in this paper are usual anti-spoofing methods (for detecting impersonation attacks), obfuscation methods are specifically discussed in Section 5: "Discussion".

The objective of this paper is to give a review of the impersonation PAD (anti-spoofing) methods that do not require any specific hardware. In other words, we focus on methods that can be implemented with only RGB cameras from Generic Consumer Devices (GCDs). This obviously raises some difficulties and limitations, *e.g.* when it comes to distinguishing between 2D planar surfaces (photo, screen) and 3D face surfaces. In the next section, we discuss the motivation for reviewing face anti-spoofing methods using only GCDs.

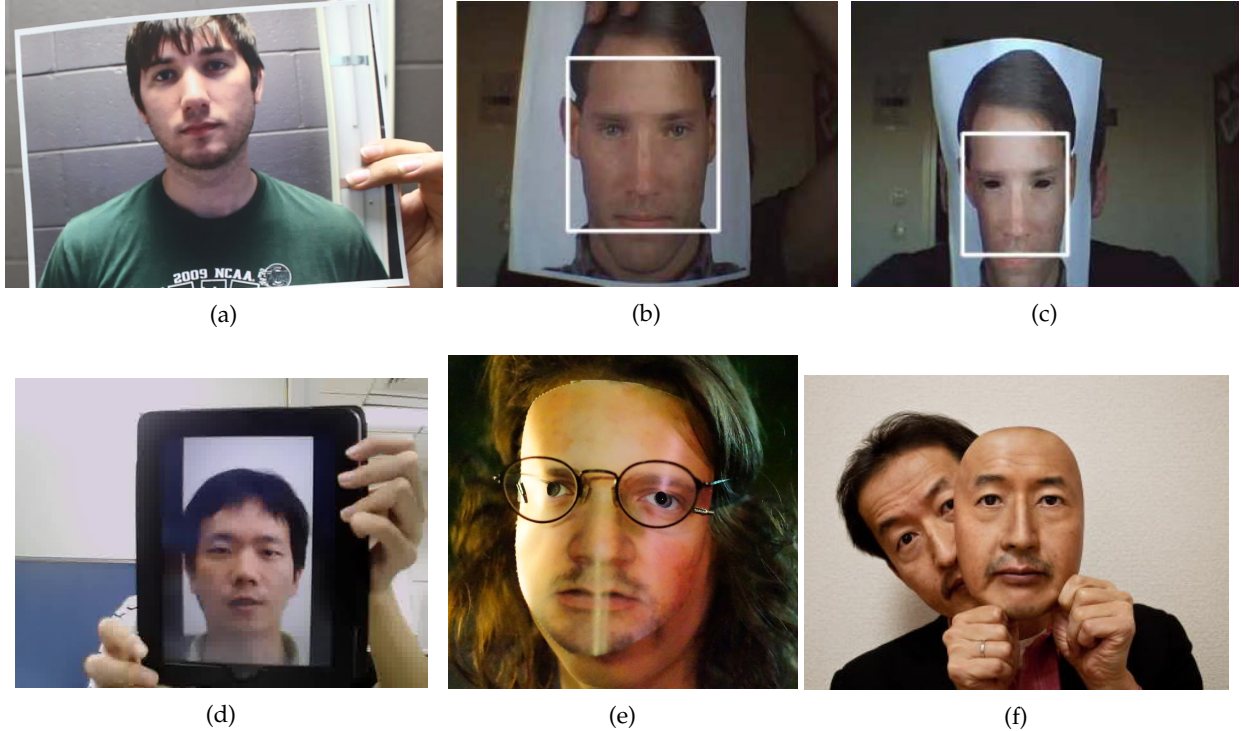


Figure 2: Examples of common face presentation attacks: (a) Printed photo attack from the SiW dataset [17]; (b) Example of a warped photo attack extracted from [18]; (c) Example of a cut photo attack extracted from [18]; (d) Video replay attack from the CASIA-FASD dataset [19]; (e) Paper-crafted mask from UMRE [20]; (f) High-quality 3D mask attack from REAL-f [17].

1.3 Face PAD methods with Generic Consumer Devices (GCD)

To the best of our knowledge, there is still no agreed-upon PAD method that can tackle all types of attacks. Given the variety of possible PAs, many face PAD approaches have been proposed in the past two decades. From a very general perspective, one can distinguish between the methods for face PAD based on specific hardware/sensors, and the approaches using only RGB cameras from GCDs.

Face PAD methods using specific hardware may rely on structured-light 3D sensors, Time of Flight (ToF) sensors, Near-infrared (NIR) sensors, thermal sensors, etc. In general, such specific sensors considerably facilitate face PAD. For instance, 3D sensors can discriminate between the 3D face and 2D planar attacks by detecting depth maps [26], while NIR sensors can easily detect video replay attacks (as electronic displays appear almost uniformly dark under NIR illumination) [27, 28, 29, 30], and thermal sensors can detect the characteristic temperature distribution for living faces [31]. Even though such approaches tend to achieve higher performance, they are not yet broadly available to the general public. Indeed, such sensors are still expensive, and rarely embedded on ordinary GCDs, with the exception of some costly devices. Therefore, the use of such specific sensors is limited to some applicative scenarios, such as physical access control to protected premises.

However, for most applicative scenarios, the user needs to be authenticated using her own device. In such scenarios, PAD methods that rely on specific hardware are therefore not usable. Thus, researchers and developers widely opt for methods based on RGB cameras that are embedded in most electronic GCDs (such as smartphones, tablets or laptops) [32, 33, 34, 35, 36, 37, 38].

This is the main reason why, in this work, we focus on the face PAD approaches that do not require any specific hardware. More precisely, we present a comprehensive review of the research work in face anti-spoofing methods for face recognition systems using only the RGB cameras of GCDs. The major contributions of this paper are listed in the section below.

1.4 Main contributions of this paper

The major contributions of this survey paper are the following:

- We propose a typology of existing face PAD methods, based on the type of PAs they aim to detect, and some specificities of the applicative scenario.
- We provide a comprehensive review of over 50 recent face PAD methods that only require (as input) images captured by RGB cameras embedded in most GCDs.
- We provide a summarized overview of the available public databases for both 2D attacks and 3D mask attacks, which are of vital importance for both model training and testing.
- We report extensive experimental results and quantitatively compare the different PAD methods under uniform benchmarks, metrics and protocols.
- We discuss some less-studied topics in the field of face PAD, such as unknown PAs and obfuscation attacks, and we provide some insights for future work.

1.5 Structure of this paper

The remainder of this paper is structured as follows. In Section 2, we propose a typology for face PAD methods based on RGB cameras from GCDs and review the most representative/recent approaches for each category. In Section 3, we present a summarized overview of the most used/interesting datasets together with their main advantages and limitations. Then, Section 4 presents a comparative experimental evaluation of the reviewed PAD methods. Section 5 provides a discussion about current trends, and some insights for future directions of research. Finally, we draw the conclusions in Section 6.

2 Overview of face PAD Methods using only RGB cameras from GCDs

2.1 Typology of face PAD methods

A variety of different typologies could be found in literature. For instance, Chingovska *et al.* [39] proposed to group the face PAD methods into three categories: motion-based, texture-based, and image-quality based methods, while Costa-Pazo *et al.* [40] considered image quality-based face PAD methods as a subclass of texture-based methods. Ramachandra and Busch [41] classified face PAD methods into two more general categories: hardware-based and software-based methods. The different approaches are then hierarchically classified into sub-classes of these two broad categories. Hernandez-Ortega *et al.* [42] divided the PAD methods as static or dynamic methods, depending on whether they take into account temporal information, or not.

Based on the type of attacks presented in section 1.2 and inspired by [41], we categorize face PAD methods into two broad categories: RGB camera-based PAD methods and PAD methods using specific hardware. As stated earlier, in this paper, we focus on the face PAD approaches that use only RGB cameras embedded in most GCDs (smartphone, tablet, laptop, *etc.*). Inside this broad category, we distinguish between the following five different classes:

1. Liveness cue-based methods;
2. Texture cue-based methods;
3. 3D geometric cue-based methods;
4. Multiple cues-based methods;
5. Methods using new trends.

As detailed in Figure 3, each of these five categories is then divided into several sub-classes, depending on the applicative scenario, or on the type of features/methods used. For each category/subcategory of PAD methods, Table 1 shows the type(s) of PAs it is aiming at detecting, whereas Figure 3 also lists all the face PAD methods that will be discussed in the remainder of this Section (over 50 methods in total).

From a very general standpoint:

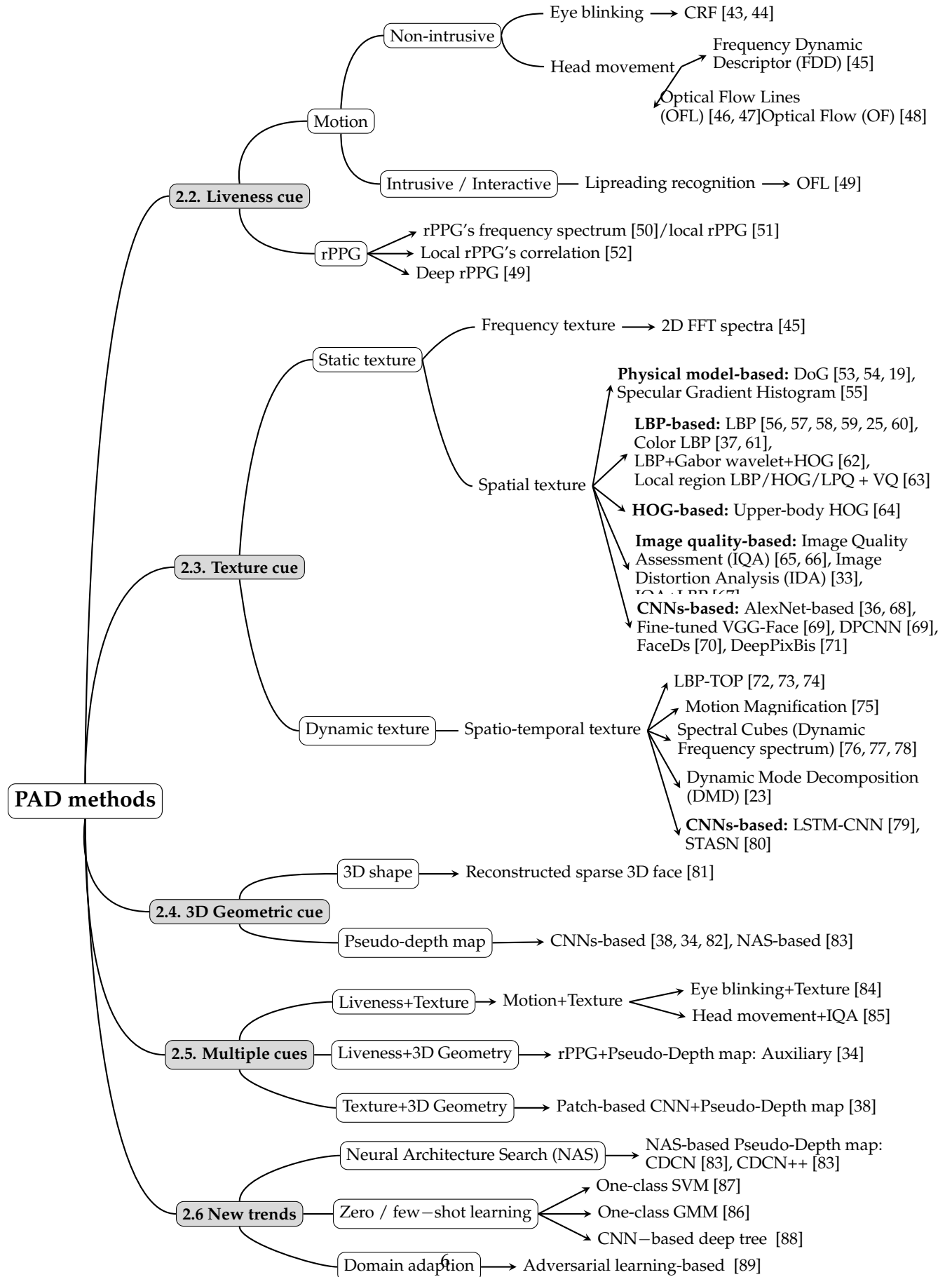


Figure 3: Our proposed typology for face PAD methods. The darkest nodes are numbered with their corresponding sub-sections in the remainder of Section 2.

Table 1: The type(s) of Presentation Attacks (PAs) each sub-type of PAD method aims to detect.

PAD methods	Sub-types	PAs
Liveness cue-based	Non-intrusive motion-based	Photo attack (except cut photo attack)
	Intrusive motion-based	Photo attack (except cut photo attack) Video replay attacks (except sophisticated DeepFakes)
	rPPG-based	Photo attack "Low quality" video replay attacks 3D mask attack (low / high quality)
Texture cue-based	Static texture-based	Photo attack Video replay attack
	Dynamic texture-based	3D mask attack (low quality)
3D Geometry cue-based	3D shape-based	Photo attack
	Pseudo-depth map-based	Video replay attack
Multiple cues-based	Liveness (Motion) + Texture	Photo attack Video replay attack
	Liveness + 3D Geometry (rPPG + Pseudo-depth map)	Photo attack Video replay attack 3D mask attack (low / high quality)
	Texture + 3D Geometry (Patched-base texture + Pseudo-depth map)	Photo attack Video replay attack

1. Liveness cue-based methods aim at detecting liveness cues in the face presentation or PAI. The most widely used liveness cues so far are motion (head movements, facial expressions, etc.) and micro intensity changes corresponding to blood pulse. Thus, liveness cue-based methods can be classified into the following two sub-categories:
 - Motion cue-based methods employ motion cues in video clips to discriminate between genuine (alive) faces and static photo attacks. Such methods can effectively detect static photo attacks, but not video replay with motion/liveness cues, and 3D mask attacks;
 - Remote PhotoPlethysmoGraphy (rPPG) is the most widely used technique for measuring face's micro intensity changes corresponding to blood pulse. rPPG cue-based methods can detect photo and 3D mask attacks, as these PAs do not show the periodic intensity changes that are characteristic of facial skin. They can also detect "low-quality" video replay attacks that are not able to display those subtle changes (due to the capture conditions and/or PAI characteristics). However, "high-quality" video replay attacks (displaying the dynamic changes of the genuine face's skin) cannot be detected by rPPG cue-based methods.
2. Texture cue-based methods use static or dynamic texture cues to detect the face PAs by analyzing the micro-texture of the surface presented to the camera. Static texture cues are generally spatial texture features that can be extracted from a single image. In contrast, dynamic texture cues usually consist in spatio-temporal texture features, extracted from an image sequence. Texture cue-based face PAD methods can detect all types of PAs. However, they might be fooled by "high-quality" 3D masks (masks with a surface texture that mimics well facial texture);
3. 3D geometric cue-based methods use 3D geometrical features, generally based on the 3D structure or depth information/map of the user's face or PAIs. 3D geometric cue-based PAD methods can detect planar photo and video replay attacks, but not (in general) 3D mask attacks;
4. Multiple cues-based methods consider different cues (e.g. motion features with texture features) to detect a wider variety of face PAs;
5. Methods using new trends do not necessarily aim at detecting specific types of PAs, but their common trait is that they rely on cutting-edge machine learning technology, such as Neural Architecture Search (NAS), zero-shot learning, domain adaption, etc.

In the remainder of this section, we present a detailed review of the over fifty recent PAD methods that are listed in Figure 3, structured using the above typology and in chronological order inside each category/sub-

category. In each category, we elaborate both the "conventional" methods that have most influenced face PAD, and their current evolutions in the deep learning era.

2.2 Liveness cue-based methods

Liveness cue-based methods are the first attempt for face PAD. Liveness cue-based methods aim to detect any dynamic physiological sign of life, such as eye blinking, mouth movement, facial expression changes and pulse beat. They can be categorized as motion-based methods (to detect the eye blinking, mouth movement and facial expression changes) and rPPG-based methods (to detect the pulse beat).

2.2.1 Motion-based methods

By detecting movements of the face/facial features, conventional motion-based methods can effectively detect static presentation attacks, such as most photo attacks (without dynamic information). However, they are generally not effective against video replay attacks that display liveness information such as eye blinking, head movements, facial expression changes, *etc.*

This is why interactive motion-based methods were later introduced, where the user is required to complete a specific (sequence of) movement(s) such as head rotation / tilting, mouth opening, *etc.* The latter methods are more effective for detecting video replay attacks, but they are intrusive for the user, unlike traditional methods that do not require the user's collaboration, and are therefore non-intrusive.

The rest of this section is structured around two sub-categories: a) non-intrusive motion-based methods, that are more user-friendly and easier to implement, and b) intrusive / interactive motion-based methods, that are more robust and can detect both static and dynamic PAs.

a) Non-intrusive motion-based methods Non-intrusive motion-based PAD methods aim at detecting intrinsic liveness based on movement (head movement, eye blinking, facial expression changes, *etc.*).

In 2004, Li *et al.* [45] first used **frequency-based** features to detect photo attacks. More specifically, they proposed the Frequency Dynamic Descriptor (FDD), based on frequency components' energy, to estimate temporal changes due to movements. By setting an FDD threshold, genuine (alive) faces can be distinguished from photo PAs, even for relatively high-resolution photo attacks. This method is easier to implement and is less computationally expensive, when compared to the previously proposed motion-based methods, that used 3D depth maps to estimate head motions [90, 91]. However, its main limitation is that it relies on the assumption that the illumination is invariant during video capture, which can't always be satisfied in a real-life scenario. This can lead to the presence of a possibly large quantity of "noise" (coming from illumination variations) in the frequency component's variations, and the method is not conceived to deal with such noise.

Unlike the approach introduced in [45], the method proposed in 2005 by Kollreider *et al.* [46, 47] works directly in the RGB representation space (and not in the frequency domain). More precisely, the authors try to detect the differences in motion patterns between 2D planar photographs and genuine (3D) faces using **optical flows**. The idea is the following: when a head has a small rotation (which is natural and unintentional), for a real face, the face parts that are nearer to the camera ("inner parts", *e.g.*, nose), move differently from the parts further away from the camera ("outer parts", *e.g.*, ears). In contrast, a translated photograph generates constant motion among all face regions [46].

More precisely, the authors proposed Optical Flow Lines (OFL), inspired from [92], to measure face motion in horizontal and vertical directions. As illustrated in Figure 4, the different greyscales obtained in the OFL from a genuine (alive face) with a subtle face rotation reflect the motion differences in between different face parts, whereas the OFL of a translated photo show constant motion.

A liveness score in $[0,1]$ is then calculated from the OFLs of the different face regions, where 1 indicates that the movement of the surface presented to the camera is coherent with a face's movement, and 0 indicates that this movement is not coherent with a face's movement. By thresholding this liveness score, the method proposed in [46] can detect printed photo attacks, even if the photo is bent, or even warped around a cylinder (as it is still far from the real 3D structure of a face). However, this method fails for most video replay attacks, and it can be disrupted by eyeglasses (because they partly cover outer parts of the face, but

are close to the camera) [47].



Figure 4: The Optical Flow Lines (OFL) images obtained from (a) a genuine (alive face) presentation with a subtle face rotation and (b) a printed photo attack with horizontal translation [47]. In (a), the inner parts are brighter than the outer parts of the face, which is characteristic of the motion differences between different face parts. In contrast, all parts of the planar photo in (b) display constant motion.

In 2009, Bao *et al.* [48] also leveraged optical flow to distinguish between 3D faces and planar photo attacks. Let us call O the object presented to the system (face or planar photo). By comparing the optical flow field of O deduced from its perspective projection to a predefined 2D object's reference optical flow field, the proposed method can determine if the given object is a really a 3D face, or a planar photo.

But, like all other optical flow-based methods, the methods in [46, 47, 48] are not robust toward background and illumination changes.

In 2007, Pan *et al.* [43, 44] chose to focus on **eye blinking**, in order to distinguish between a face and a facial photo. Eye blinking is a physiological behavior that normally happens 15 to 30 times per minute [93]. Therefore, it is possible for GCD cameras having at least 15 frames per second (fps) – which is almost all GCD cameras – to capture two or more frames per blink [43]. Pan *et al.* [44, 43] proposed to use Conditional Random Fields (CRFs) fed by temporally observed images x_i to model eye-blinking with its different estimated (hidden) states y_i : Non-Closed (NC, including opened and half-opened eyes), then Closed (C), then NC again, as illustrated in Figure 5. The authors showed that their CRF-based model (discriminative model) is more effective than Hidden Markov Model (HMM) based methods [94] (generative model), as it takes into account longer-range dependencies in the data sequence. They also showed that their model is superior to another discriminative model: Viola and Jones' Adaboost cascade [95] (as the latter is not conceived for sequential data). This method, like all other methods based on eye-blinking, can effectively detect printed photo attacks, but not video replay attacks or eye-cut photo attacks that simulate blinking [19].

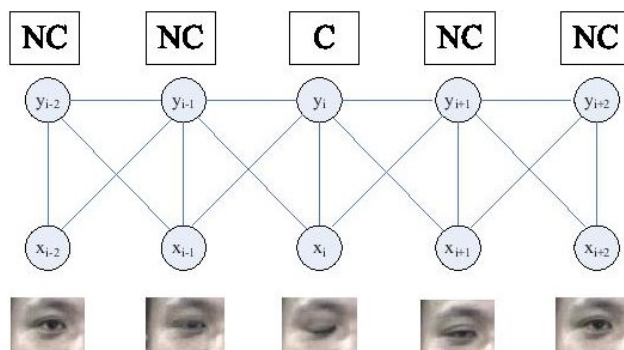


Figure 5: CRFs-based eye-blinking model [44, 43]. In this example each hidden state y_i is conditioned by its corresponding observation (image) x_i and its two neighboring observations x_{i-1} and x_{i+1} .

b) Intrusive motion-based methods Intrusive methods (also called interactive methods) are usually based on a Challenge-Response mechanism that requires the users to satisfy the system's requirements. In this paragraph, we present methods where the challenge is based on some pre-defined head/face movement (*e.g.*, blinking the eyes, moving the head in a certain way, adopting a given facial expression or uttering a

certain sequence of words).

In 2007, Kollreider *et al.* [49] first proposed an interactive method that can detect replay attacks as well as photo attacks by reading on the presented face's lips when the user is prompted to utter a randomly-determined sequence of digits. Like in their previous work [46] mentioned above, Kollreider *et al.* [49] use OFL to extract mouth motion. An interesting feature of this method is that it combines face detection and face PAD in a holistic system. Thus, the integrated face detection module can also be used to detect the mouth region, and OFL is used for both face detection and face PAD.

Then, a 10-class Support Vector Machines (SVM) [96] is trained from 160-dimensional velocity vectors extracted from the mouth region's OFL so as to perform recognition of the 0-9 digital. This method detects effectively printed photo attacks and most video replay attacks. However, it is vulnerable to mouth-cut photo attacks and, even though this topic was not studied yet (at least, to the best of our knowledge), it certainly cannot detect sophisticated DeepFakes [97, 98, 99] where the impostor can "play" on-demand any digit. Another limitation of this method is that it is based on visual cues only, *i.e.* it does not consider audio together with images (unlike multi-modal audio-visual methods [100]). This makes it vulnerable to "visual-only" DeepFakes, which are of course easier to obtain than realistic audio-visual DeepFakes (with both the facial features and voice of the impersonated genuine user).

More generally, the emerging technology of DeepFakes [101] is a great challenge for interactive motion-based PAD methods. Indeed, based on deep learning models such as autoencoders and generative adversarial networks [102, 103], DeepFakes can superimpose face images of a target person to a video of a source person in order to create a video of the target person doing or saying things that the source person does or says [97, 98, 99]. Impostors can therefore use DeepFake generation apps like FaceApp [104] or FakeApp [105] to easily create a video replay attack showing a genuine user's face satisfying the system requirements during the Challenge-Response authentication. Interactive motion-based methods generally have difficulties to detect DeepFakes-based video replay attacks.

However, recent works show that rPPG-based [106, 107] and texture-based methods [108] can be used to detect video attacks generated using DeepFakes. rPPG-based methods are discussed in the next section.

2.2.2 Liveness detection based on Remote PhotoPlethysmoGraphy (rPPG)

Unlike the head/face movements that are relatively easy to detect, the intensity changes in the facial skin that are characteristic of pulse/heartbeat are imperceptible for most human eyes. So as to detect automatically these subtle changes, remote PhotoPlethysmoGraphy (rPPG) was proposed [50, 17, 34]. rPPG can detect blood flow using only RGB images from a distance (in a non-intrusive way), based on the analysis of the variations in the absorption and reflection of light passing through human skin. The idea behind rPPG is illustrated in Figure 6.

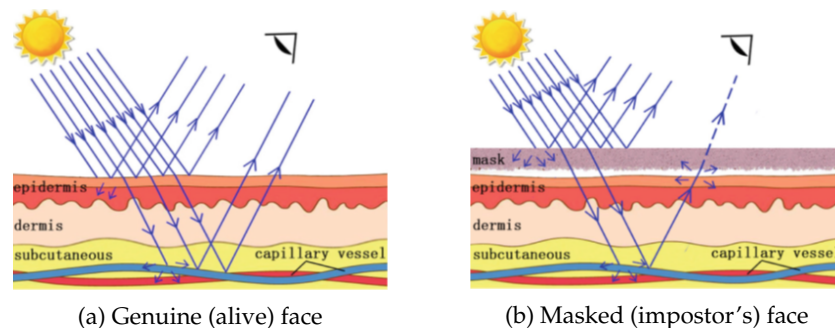


Figure 6: Illustration of how rPPG can be used to detect blood flow for face PAD [17]. (a) On a genuine (alive) face, the light penetrates the skin and illuminates capillary vessels in the subcutaneous layer. Blood oxygen saturation changes within each cardiac cycle, leading to periodic variations in the skin's absorption and reflection of the light. These variations are observable by RGB cameras. (b) On a masked face, the mask's material blocks the light absorption and reflection, leading to no (or insignificant) variations in the reflected light.

Since photo-based PAs do not display any periodical variation in the rPPG signal, they can be detected easily by rPPG-based methods. Moreover, as illustrated in Figure 6, most kinds of 3D masks (including high-quality masks, except maybe extremely thin masks) can be detected by rPPG-based methods. However, "high-quality" video replay attacks (with good capture conditions and good-quality PAI) can also display the periodic variation of the genuine face's skin light absorption/reflection. Thus, rPPG-based methods are only capable of detecting low-quality video replay attacks.

The first methods that applied rPPG to face PAD were published in 2016. Li *et al.* [50] proposed a simple approach whose framework is shown in Figure 7. The lower half of the face is detected and extracted as a Region of Interest (RoI). The rPPG signal is composed by the average RGB value of pixels in the ROI for each RGB channel of each video frame. This rPPG signal is then filtered (to remove noise and to extract the normal pulse range) and transformed into a frequency signal by Fast Fourier Transform (FFT). Two frequency features per channel (denoted as $[E_r, E_g, E_b]$ and $[\Gamma_r, \Gamma_g, \Gamma_b]$ in Figure 7) are extracted for each color channel based on the Power Spectral Density (PSD). Finally, these (concatenated) feature vectors are fed into a SVM to differentiate the genuine face presentation from PAs.

This rPPG-based method can effectively detect photo-based and 3D mask attacks – even high-quality 3D masks – but not (in general) video replay attacks. Because, on the other hand, texture-based methods (see Section 2.3) can detect video replay attacks but not realistic 3D masks [56, 25, 37], the authors also proposed a cascade system that uses first their rPPG-based method (to filter photo or 3D mask attacks) and then a texture-based method (to detect video replay attacks).

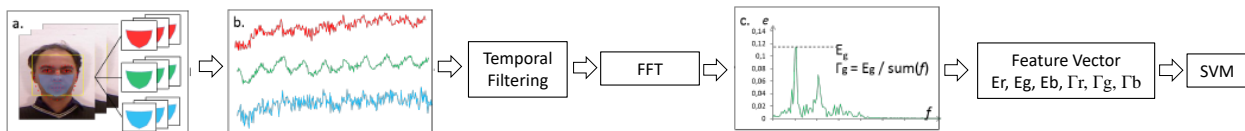


Figure 7: Framework of the rPPG-based method proposed in [50] (Figure extracted from [50]).

Also in 2016, Liu *et al.* [52] proposed another rPPG-based method for detecting 3D mask attacks. Its principle is illustrated in Figure 8. This method has three interesting features compared to the above-mentioned approach proposed in [50] the same year. First, rPPG signals were extracted from multiple facial regions instead of just the lower half of the face. Secondly, the correlation of any two local rPPG signals was used as a discriminative feature (assuming they should all be consistent with the heartbeat's rhythm). Thirdly, a confidence map is learned, so as to weight each region's contribution: robust regions that contain strong heartbeat signals are emphasized, whereas unreliable regions containing less heartbeat signals (or more noise) are weakened.

Finally, the weighted local correlation-based features are fed into a SVM (with RBF kernel) to detect photo and 3D mask PAs. This approach is more effective than the one proposed by Li *et al.* in [50].

In 2017, Nowara *et al.* [51] proposed PPGSecure, a local rPPG-based approach within a framework that is very similar to the one in [50] (see Figure 7). The rPPG signals are extracted from three facial regions (forehead and left/right cheeks) and two background regions (on the left and right of the facial region). The use of background regions provides robustness against noise due to illumination fluctuations, as this noise can be subtracted from the facial regions after having been detected in the background regions. Finally, the Fourier spectrum's magnitudes of the filtered rPPG signals are fed into a SVM or a Random Forest Classifier [109]. The authors showed experimentally the interest of using background regions, and their method obtained better performances than the one in [52] on some dataset.

In 2018, Liu *et al.* [34] proposed a deep learning-based approach that can learn rPPG signals in a robust way (under different poses, illumination conditions and facial expression). In this approach, rPPG estimations (pseudo-rPPGs) were combined with the estimations of 3D geometric cues, in order to tackle not only photo and 3D mask attacks (like all rPPG-based methods), but also video replay attacks. Therefore, this approach is detailed together with other multiple cue-based approaches, in Section 2.5.2, on page 23.

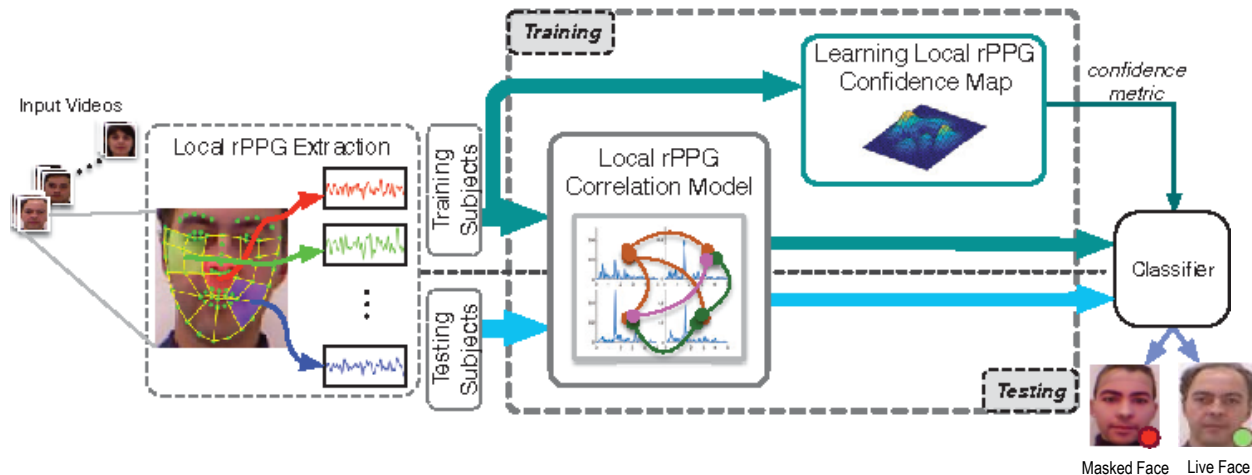


Figure 8: Framework of the local rPPG correlation-based method proposed in [52] (Figure extracted from [52]).

As mentioned in the previous section, recent studies show that, unlike motion-based PAD methods, rPPG-based PAD methods can be used to detect DeepFake videos.

Indeed, in 2019, Fernandes *et al.* [106] proposed to use Neural Ordinary Differential Equations (Neural-ODE) [110] for heart rate prediction. The model is trained on the heart rate extracted from the original videos. Then the trained Neural-ODE is used to predict the heart rate of Deepfake videos generated from these original videos. The authors show that there is a significant difference between the original videos' heart rates and their predictions in the case of DeepFakes, implicitly showing that their method could discriminate between Deepfakes and genuine videos.

A more sophisticated method was proposed in 2020 by Ciftci *et al.* [107], in which several biological signals (including the rPPG signal) are fed into a specifically designed Convolutional Neural Network (CNN) to discriminate between genuine videos and DeepFakes. The reported experimental results are very encouraging.

2.3 Texture cue-based methods

Texture feature-based methods are the most widely used for face PAD so far. Indeed, they have several advantages compared to other kinds of methods. First, they are inherently non-intrusive. Second, they are capable to detect almost any kind of known attacks, *e.g.* photo-based attacks, video replay attacks and even some 3D mask attacks.

Unlike the liveness cue-based methods that rely on dynamic physiological signs of life, texture cue-based methods explore the texture properties of the object presented to the system (genuine face or PAI). With texture cue-based methods, PAD is usually formalized as a binary classification problem (real face / non face) and these methods generally rely on a discriminative model.

Texture cue-based methods can be categorized as static texture-based, and dynamic texture-based. Static texture-based methods extract spatial or frequential features, generally from a single image. In contrast, dynamic texture-based methods explore spatio-temporal features extracted from video sequences. The next two sub-sections present the most prominent approaches from these two types.

2.3.1 Static texture-based methods

The first attempt to use static texture clues for face PAD dates back to 2004 [45]. In this method, the difference of light reflectivity between a genuine (alive) face and its printed photo is analyzed using their **frequency representations** (and, more specifically, their 2D Fourier spectra). Indeed, as illustrated in Figure 9, the 2D

Fourier spectrum of a face picture has much less high-frequency components than the 2D Fourier spectrum of a genuine (alive) face image.

More specifically, the method relies on High-Frequency Descriptors (HFD), defined as the energy percentage explained by high-frequency components in the 2D Fourier spectrum. Then, printed photo attacks are detected by thresholding the HFD value (attacks being below the threshold). This method works well only for small images with poor resolution. For instance, it is vulnerable to photos of 124x84mm or with 600 dpi resolution.



Figure 9: From left to right: : a genuine (alive) face, a printed photo attack, and their respective 2D Fourier spectra (Figure extracted from [45]).

In 2010, Tan *et al.* [53] first modeled the respective reflectivities of images of genuine (alive) faces and face printed photos using **physical models** (here Lambertian models) [111], in which latent samples are derived using Difference of Gaussian (DoG) filtering [112]. The idea behind this method is that an image of a face printed photo tends to be more distorted than an image of real face, because it has been captured twice (by possibly different sensors) and printed once in the meanwhile (see Section 3.1 for more information about the capture process), whereas real faces are only captured once (by the biometric system only). Several classifiers were tested, among which Sparse Nonlinear Logistic Regression (SNLR) and SVMs, with SNLR proving to be slightly more effective.

Since DoG filtering is sensitive to illumination variations and partial occlusion, Peixoto *et al.* [54] proposed in 2011 to apply Contrast-Limited Adaptive Histogram Equalization (CLAHE) [113] to pre-process all images, showing the superiority of CLAHE to a simple histogram equalization.

Similar to the 2010's work from Tan *et al.*, Bai *et al.* [55] also used, the same year, a physical model to analyze the images' micro-textures, using Bidirectional Reflectance Distribution Functions (BRDF). The original image's normalized specular component (called *specular ratio image*) is extracted [114], then its gradient histogram (called *specular gradient histogram*) is calculated. As shown in Figure 10, the shapes of the specular gradient histograms of a genuine (alive) face and of a printed photo are quite different. To characterize the shape of a specular gradient histogram, a Rayleigh histogram model is fitted on the gradient histogram. Then, its two estimated parameters σ and β are used to feed a SVM. This SVM is trained to discriminate between genuine face images and planar PAs (in particular, printed photos and video replay attacks).

As shown in Figure 10, this method can detect planar attacks just from a small patch of the image. But, specular component extraction requires a highly contrasted image, and therefore this method is vulnerable towards any kind of blur.

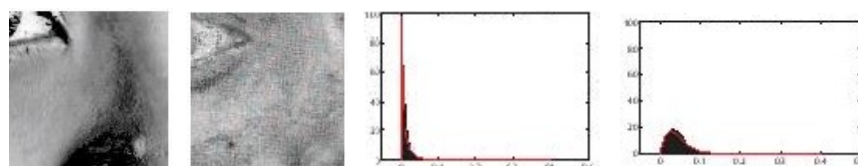


Figure 10: From left to right: patches extracted from a genuine face image and a printed photo, and their respective specular gradient histograms.

Local Binary Pattern (LBP) [115] is one of the most widely used hand-crafted texture features in the face analysis-related problems, such as face recognition [116], face detection [117] and facial expression

recognition [118]. Indeed, it has several advantages, including a certain robustness toward illumination variations.

In 2011, Määttä *et al.* [56] first proposed to apply multi-scale LBP to face PAD. Unlike the previously described static texture-based approaches [53, 55], LBP-based methods do not rely on any physical model; they just assume that the differences in surface properties and light reflection between a genuine face and a planar attack can be captured by the LBP features.

Figure 11 illustrates this method. Three different LBPs were applied on a normalized 64x64 image in [56]: $LBP_{8,2}^{u2}$, a uniform circular LBP extracted from 8 pixel-neighbourhood with 2-pixel radius, $LBP_{16,2}^{u2}$, a uniform circular LBP extracted from 16 pixel-neighbourhood with 2-pixel radius, and $LBP_{8,1}^{u2}$, a uniform circular LBP extracted from 8 pixel-neighbourhood with 1-pixel radius. Finally, a concatenation of all generated histograms formed a 833-bin/dimension histogram. This histogram is then used as a global micro-texture feature, and fed to a non-linear (RBF) SVM classifier for face PAD.

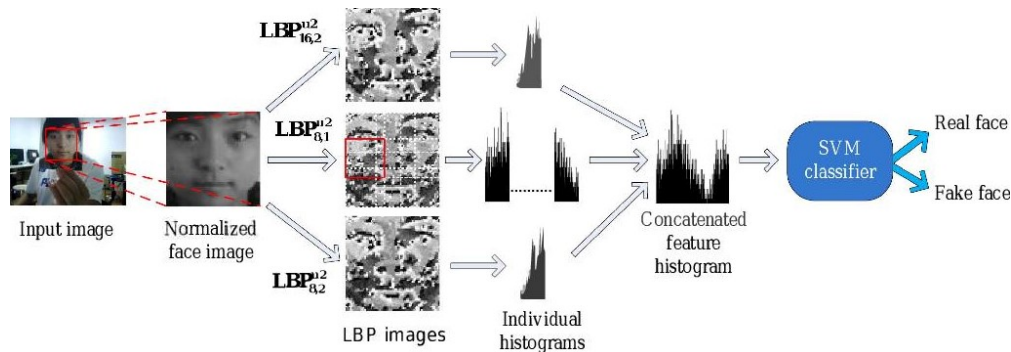


Figure 11: Illustration of the approach proposed in [56]. Firstly, the face is detected, cropped and normalized into a 64×64 pixel image. Then, $LBP_{8,2}^{u2}$ and $LBP_{16,2}^{u2}$ are applied on the normalized face image, which generates a 59-bin histogram and a 243-bin histogram respectively. The obtained $LBP_{8,1}^{u2}$ image is also divided into 3×3 overlapping regions (as shown in the middle row). As each region generates a 59-bin histogram, a single 531-bin histogram is obtained by their concatenation. Then, all individual histograms are concatenated to obtain a 833-bin/dimension ($59+243+531$) histogram, which is fed to a nonlinear SVM classifier, to detect photo / video replay attacks.

In 2012, the authors extended their work in [62], adding two more texture features within the same framework: Gabor wavelets [119] (that can describe facial macroscopic information), and Histogram of Oriented Gradients (HOG) [120] (that can capture the face’s edges or gradient structures). Each feature (LBP-based global micro-texture feature, Gabor wavelets and HOG) is transformed into a compact linear representation by using a homogeneous kernel map function [121]. Then each transformed feature is separately fed into a fast linear SVM. Finally, late fusion between the scores of the three SVM output is applied, so as to generate a final decision. The authors showed the superiority of this approach compared to the method they previously introduced in [56].

In 2013, the same authors continued to extend their work [64], using this time the upper-body region instead of the face region to detect spoofing attacks. As shown in Figure 12, the upper-body region includes more scenic cues of the context, which enables to detect the boundaries of the PAI (*e.g.*, video screen frame or photograph edge), and, possibly, the impostor’s hand(s) holding the PAI. As a local shape feature, HOG is calculated from the upper-body region to capture the continuous edges of the PAI (see Figure 12(d)). Then, this HOG feature is fed to a linear SVM for detecting photo or video replay attacks. The upper-body region is detected using the method in [122], that can also be used to filter poor attacks (where the PAI is poorly positioned or with strong discontinuities between the face and shoulder regions), as shown in Figure 12(c).

In the same spirit of using context surrounding the face, Yang *et al.* [63] proposed (also in 2013) to use a $1.6 \times$ enlarged face region, called Holistic-Face (H-Face), to perform PAD. In order to focus on the facial regions that play the most important role in face PAD, the authors segmented four canonical facial regions: the left eye region, right eye region, nose region, and mouth region, as shown in Figure 13. The rest of the face (mainly the facial contour region) and the original enlarged face images were divided as 2×2

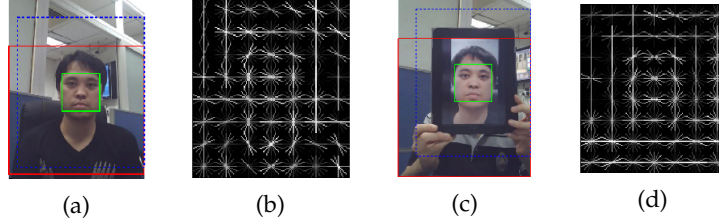


Figure 12: Examples of upper-body images from CASIA-FASD [19] and their HOG features [64]: (a) Upper-body of a genuine face; (b) HOG feature of the blue dashed rectangle in (a); (c) Video replay attack; (d) HOG feature of the blue dashed rectangle in (c). Figure extracted from [64].

blocks respectively, to obtain another eight components. Thus, twelve face components are used in total. Then, different texture features such as LBP [56], HOG [120] and Local Phase Quantization (LPQ) [123] are extracted as low-level features from each component. Instead of directly feeding the low-level features to the classifier, a high-level descriptor is generated based on the low-level features by using spatial pyramids [124] with a 512-word codebook. Then the high-level descriptors are weighted using average pooling to extract higher-level image representations. Finally, the histogram of these image representations are concatenated into a single feature vector fed into a SVM classifier to detect PAs.

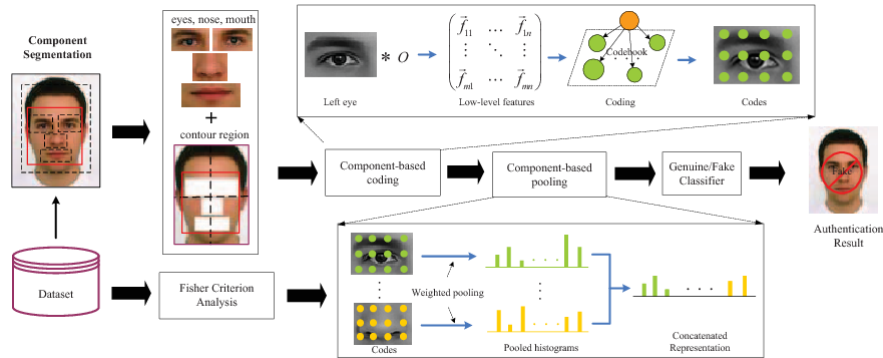


Figure 13: Illustration of the approach proposed in [63] (Figure extracted from [63]).

In 2013, Kose *et al.* [57] first proposed a static texture-based approach to detect 3D mask attacks. Due to the unavailability of public mask attacks databases at that time, 3D mask PADs were much less studied than photo or video replay attacks. In this work, the LBP-based method in [56] is directly applied to detect 3D mask attacks by using the texture (original) image or depth image of the 3D mask attacks (from a self-constructed database), as shown in Figure 14. Note that all the texture images and the depth images were obtained by MORPHO². This work showed that using the texture (original) image is better than using a depth image for detecting 3D mask attacks with LBP features. The authors also proposed in [58] to improve this method by fusing the LBP features of the texture image and depth image. They showed the superiority of this approach toward the previous method (using only texture images).

Erdogmus *et al.* [59, 25] also proposed in 2013 a method for 3D mask attack detection based on LBP. They used different classifiers, such as Linear Discriminant Analysis (LDA) and SVMs. On the proposed 3D Mask Attack Database (3DMAD), which is also the first public spoofing database for 3D mask attacks, LDA was proved to be best among the tested classifiers.

Galbally *et al.* [65, 66] introduced in 2013 and 2014 new face PAD methods based on **Image Quality Assessment (IQA)**, assuming that a spoofing image captured in a photo or video replay PA should have a different quality than a real sample, as it was captured twice instead of once for genuine faces (this idea is similar to the underlying idea of the method in [53] presented above). The quality differences concern

²<http://www.morpho.com/>



Figure 14: The texture image and its corresponding depth image, for (a) a real access and (b) 3D mask attack. Figure extracted from [25].

sharpness, colour and luminance levels, structural distortions, etc. 14 image quality measures and 25 image quality measures were adopted in [65] and [66] respectively, to assess the image quality using scores extracted from single images. Then, the image quality scores were combined as a single feature vector and fed into a LD or Quadratic Discriminant Analysis (QDA) classifier, to perform face PAD. The major advantage of the IQA-based methods is that it is not a trait-specific method, *i.e.* it does not rely on priori face/body region detection, so this is a "multi-biometric" method that can also be employed for iris or fingerprint-based liveness detection. However, the performance of the proposed IQA-based methods for PAD was limited compared to other texture-based methods, and the method is not conceived to detect 3D mask attacks.

In 2015, Wen *et al.* [33] also proposed an IQA-based method, using the analysis of image distortion, for face PAD. Unlike the methods from Tan *et al.* [112] and Bai *et al.* [55] presented above – methods that work in the RGB space – this method analyzes the image chromaticity and the colour diversity distortion in the HSV (Hue, Saturation, and Value) space. Indeed, when the input image resolution is not enough, it is hard to tell the difference between a genuine face and a PA based only on the RGB image (or grey-scale image). The idea here is to detect imperfect/limited colour rendering of a printer or LCD screen. A 121-dimensional image distortion feature (which consists of a three-dimension specular reflection feature [125], two-dimension no-reference blurriness feature [126, 127], a 15-dimension chromatic moment feature [128] and a 101-dimension colour diversity feature) is fed into two SVMs corresponding respectively to photo attacks and video replay attacks. Finally, a score-level fusion based on the Min Rule [129] gives the final result. Unlike the IQA score-based features used in [66], this feature is face-specific. The proposed method has shown a promising generalization performance, when compared with other texture-based PAD methods.

In 2015, Boulkenafet *et al.* [37, 61] also proposed to extract LBP features in HSV or YCbCr colour spaces. Indeed, subtle differences between a genuine face and a PA can be captured by chroma characteristics, such as the Cr channel that is separated from the luminance in the YCbCr colour space (see Figure 15). Only by simply changing the colour space used, this LBP-based method achieved state-of-the-art performances, when compared with some much more complicated PAD methods based on Component Dependent Descriptor (CDD) [63] and even the emerging deep CNNs [36]. This work showed the interest of using diverse colour spaces for face PAD.

In 2016, Patel *et al.* [67] first proposed a spoof detection approach on a smartphone. They used a concatenation of multi-scale LBP [56] and image quality-based colour moment features [33] as a single feature vector fed into a SVM for face PAD. Like in [64], this work also introduced a strategy to pre-filter the poor attacks before employing the sophisticated SVM for face PAD. For this purpose, the authors proposed to detect the bezel of PAI (*e.g.*, the white bezel of photo or the screen black bezel along the border) and the Inter-Pupillary Distance (IPD). For bezel detection, if the pixel intensity values remain fairly consistent (over 60 or 50 pixels) on any four sides (top, bottom, left, and right sides), the region is considered as belonging to the bezel of a PAI. For IPD detection, if the IPD is too small (*i.e.* the PAI is too far from the acquisition camera) or too large (*i.e.* the PAI is too close to the acquisition camera), then the presentation is classified as a PA. The threshold is set to a difference exceeding two times the IPD's standard deviation observed for genuine faces with a smartphone. This strategy, relying on two simple counter-measures, can efficiently filter almost 95% of the poorest attacks. However, it may generate false rejections (*e.g.* if the genuine user is

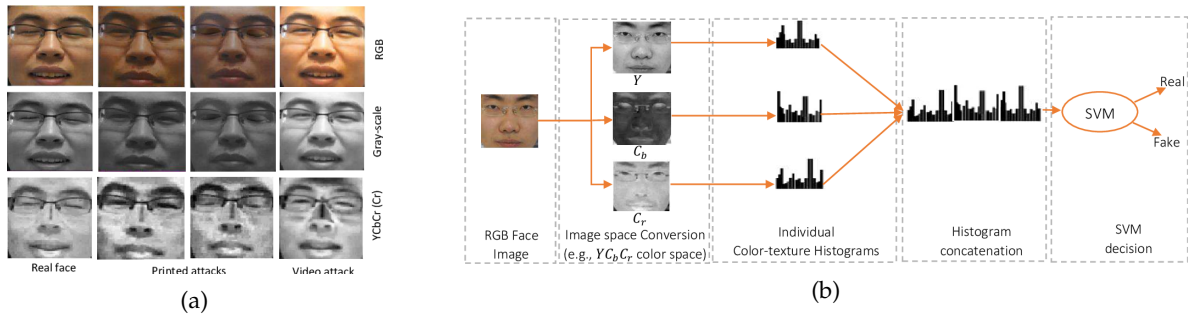


Figure 15: Illustration of the method in [61]. (a) Example of a genuine face presentation, and its corresponding printed photo and video attacks in RGB, greyscale and YCbCr colour space. (b) Architecture of the method proposed in [61].

wearing a black t-shirt on a dark background).

More recently, **deep learning-based methods** are used to learn automatically the texture features. Researchers studying these techniques rather focus on designing an appropriate neural network so as to learn the best texture features, than to design the texture features themselves (as it is the case with most hand-crafted features presented above).

The first attempt to use Convolutional Neural Networks (CNNs) for detecting spoofing attacks was claimed by Yang *et al.* for their 2014 method [36]. In this method, a one-path AlexNet [2] is used for learning the texture features that best discriminate PAs (see Figure 16). The usual output of AlexNet (a 1000-way softmax) is replaced by a SVM with binary classes. The fully-connected bottleneck layer, *i.e.* fc7, is extracted as the learned texture feature and fed into the binary SVM. Instead of being an end-to-end framework like many CNN frameworks used nowadays, the proposed approach was basically using a quite conventional SVM-based general framework, only replacing the hand-crafted features by the features learned by AlexNet. This method was shown to attain significant improvements when the input image was enlarged by a scale of 1.8 or 2.6. These results are consistent with the previous studies in [64, 63], that had already shown that including more context information from the background can help face PAD. It was the first time that CNNs were proven to be effective for automatically learning texture features for face PAD. This method has surpassed almost all the existing state-of-the-art methods for photo and video replay attacks. It showed the potential of deep CNNs for face PAD. Later, more and more CNN-based methods were explored for face PAD.

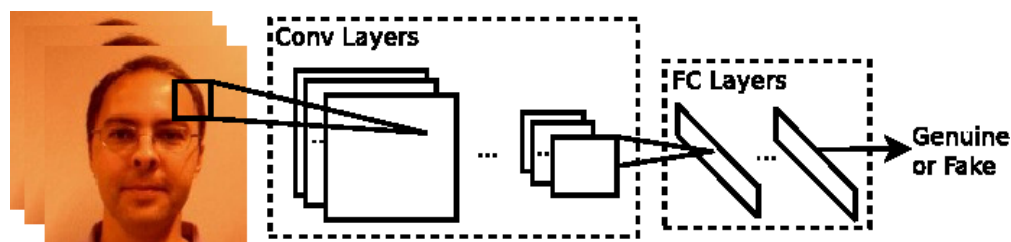


Figure 16: Illustration of the method proposed in [36] (Figure extracted from [36]).

In 2016, Patel *et al.* [68] first proposed an end-to-end framework based on one-path AlexNet [2], namely CaffeNet, for face PAD. A two-way softmax replaced the original 1000-way softmax as a binary classifier. Given the small sizes of existing face spoof databases, especially at that time, such deep CNNs were likely to overfit [68] if trained on such datasets. Therefore, the proposed CNN was pre-trained on ImageNet [130] and WebFace [131] to provide a reasonable initialization, and fine-tuned using the existing face PAD databases. More specifically, two separate CNNs are trained, respectively from aligned face images and enlarged images including some background. Finally, a voting fusion is used to generate a final decision. Just like Yang *et al.*'s method [36], the proposed CNN-based method has surpassed the state-of-the-art methods based on hand-crafted features for photo and video replay attacks.

Also in 2016, Li *et al.* [69] proposed to train a deep CNN based on VGG-Face [7] for face PAD. As in [68], the CNN was pre-trained on massive datasets, and fine-tuned on the (way smaller) face spoofing database. Furthermore, the features extracted from the different layers of the CNN were fused to a single feature, fed into a SVM for face PAD. However, as the dimension of the fused feature is much higher than the number of training samples, this approach is prone to overfitting. Principal component analysis (PCA) and the so-called *part features* are therefore used to reduce the feature dimension. To obtain part features, the mean feature map in a given layer is firstly calculated. Then, the critical positions in the mean feature map are selected, in which the values are higher than 0.9 times the maximum value in the mean feature map. Finally, the values of the critical positions on each feature map are selected to generate the part feature. The concatenation of all part features of all feature maps is used as the global part feature. Then PCA is applied on the global part feature to further reduce the dimension. Finally, the condensed part feature is fed into a SVM to discriminate between genuine (real) faces and PAs. Benefitting from using a deeper CNN based on VGG-Face, the proposed method has achieved state-of-the-art performances, in both intra-dataset and cross-dataset scenarios (see section 4), for detecting photo and video replay attacks.

In 2018, Jourabloo *et al.* [70] proposed to estimate the noise of a given spoof face image to detect photo / video replay attacks (the authors also claimed that the proposed method could be applied to detect makeup attacks). In this work, the spoof image was regarded as the summation of the genuine image and image-dependent noise (*e.g.*, blurring, reflection and moiré pattern) introduced when generating the spoof image. Since the noise of a genuine image was assumed as zero in this work, a spoof image can be detected by thresholding the estimated noise. A GAN framework based on CNNs, De-Spoof Net (DS Net), was proposed to estimate the noise. However, as there is no noise ground-truth, instead of assessing the quality of noise estimation the authors de-noise the spoof images and assess the quality of the recovered (de-noised) image using Discriminative Quality Net (DQ Net) and Visual Quality Net (VQ Net). Besides, by fusing different losses for modelling different noise patterns in DS Net, the proposed method has shown a superior performance compared to other state-of-the-art deep face PAD method such as [34].

In 2019, George *et al.* [71] proposed Deep Pixelwise Binary Supervision (DeepPixBiS), based on DenseNet [132], for face PAD. Instead of only using the binary cross-entropy loss of the final output (denoted as Loss 2 in Figure 17) as in [68], DeepPixBiS also uses during training a pixel-wise binary cross-entropy loss based on the last feature map (denoted as Loss 1 in Figure 17). Each pixel in the feature map is annotated as 1 for a genuine face input and 0 for a spoof face input. In the evaluation / test phase, only the mean value of pixels in the feature map is used as the score for face PAD. Thanks to the powerful DenseNet and the proposed pixel-wise loss forcing the network to learn the patch-wise feature, DeepPixBiS showed a promising PAD performance for both photo and video replay attacks.

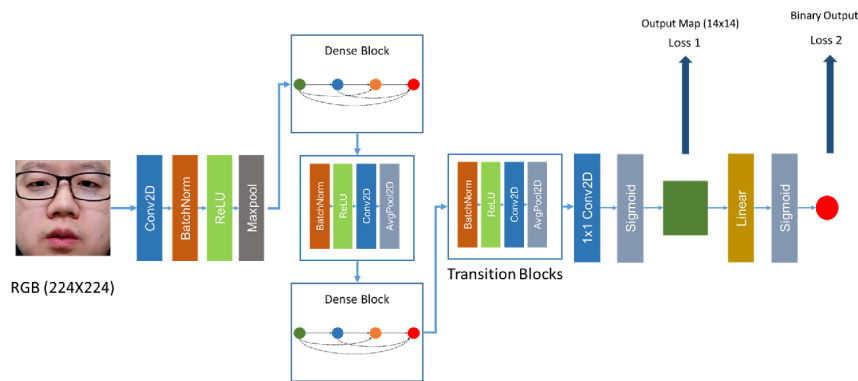


Figure 17: The diagram of Deep Pixelwise Binary Supervision (DeepPixBiS) as shown in [71].

2.3.2 Dynamic texture-based methods

Unlike the static texture-based methods, that extract spatial features usually based on a single image, dynamic texture-based methods extract spatio-temporal features using an image sequence.

Pereira *et al.* [72, 73] first proposed, in 2012 and 2014, to apply a dynamic texture based on **LBP** [115], for face PAD. More precisely, they introduced the LBP from Three Orthogonal Planes (LBP-TOP) feature [118]. LBP-TOP is a spatio-temporal texture feature extracted from an image sequence considering three orthogonal planes intersecting at the current pixel in the XY direction (as in traditional LBP), but also in XT and YT directions, where T is the time axis of the image sequence, as shown in Figure 18. The sizes of the XT and YT planes depend on the radii in the direction of time axis T, which is indeed the number of frames before or after the central frame in the image sequence. Then the conventional LBP operation can be applied to each of the three planes. The concatenation of the three LBP features extracted from the three orthogonal planes generates the LBP-TOP feature of the current image. Similarly to many static LBP-based face PAD methods, LBP-TOP is then fed into a classifier such as SVM, LDA or χ^2 distance-based classifier to perform face PAD.

In 2013, the same authors extended their work [74] by proposing two new training strategies to improve generalization. One strategy was to train the model with the combination of multiple datasets. The other was to use a score level fusion-based framework, in which the model was trained on each dataset, and a sum of the normalized score of each trained model was used as the final output. Despite the fact that these two strategies somehow ameliorate generalization, they have obvious drawbacks. First, even a combination of multiple databases cannot deal with new types of attacks that are not included in the current training datasets, so the model has to be re-trained when a new attack type is added. Second, the fusion strategy relies on an assumption of statistical independence that is not necessarily verified in practice.

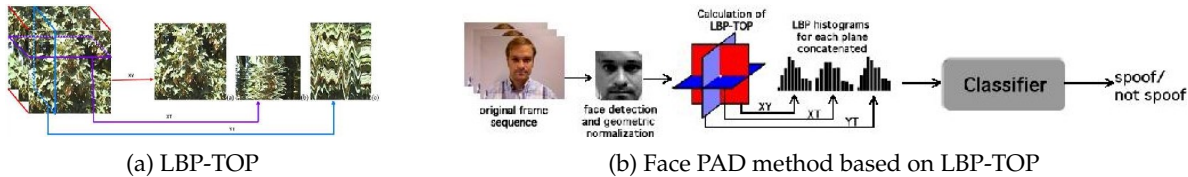


Figure 18: Illustration of the LBP-TOP from [118] and LBP-TOP based face PAD method [72]. (a) Three orthogonal planes, *i.e.* XY plane, XT plane and YT plane, of LBP-TOP features extracted from an image sequence; (b) the framework of the approach based on LBP-TOP introduced in [72].

Also in 2013, Bharadwaj *et al.* [75] proposed to use motion magnification [133] as a preprocessing, to enhance the intensity value of motion in the video before extracting the texture features. The authors claimed that the motion magnification might enrich the texture of the magnified video. The authors proposed to apply Histogram of Oriented Optical Flows (HOOF) [134] on the enhanced video to conduct face PAD. HOOF calculates the **optical flow** between frames at a fixed interval and collects the optical flow orientation angle weighted by its magnitude in a histogram. The histogram is computed from local blocks, and the resulting histogram for each block are concatenated to form a single feature vector as shown in Figure 19. HOOF is much computationally lighter than LBP-TOP. However, the proposed method based on motion magnification needs to accumulate a large number of video frames (> 200 frames), which makes it hardly applicable real-time, resulting in solutions that are not very user-friendly.

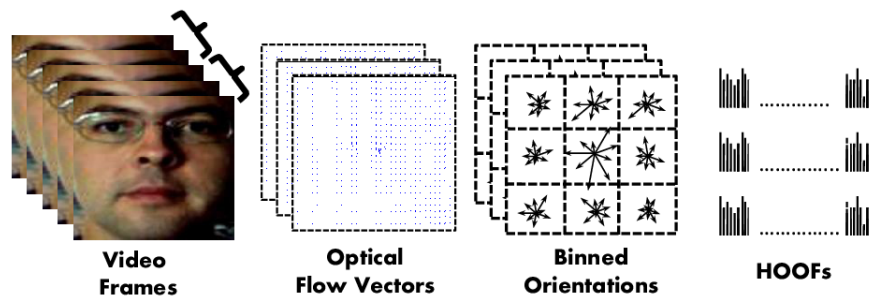


Figure 19: Illustration of the Histogram of Oriented Optical Flows (HOOF) feature proposed in [75].

In 2012 and 2015, Pinto *et al.* [76, 77, 78] proposed a PAD method based on the analysis of a video's **Fourier spectrum**. Instead of analyzing Fourier spectra of the original video as in Li *et al.* [45], the proposed method

analysed the Fourier spectra of the residual noise videos, which only include the noise information. The objective is to capture the effect of the noise introduced by the spoofing attack, *e.g.*, the moiré pattern effect shown in Figure 20(b) and Figure 20(c). In order to obtain a residual noise video, the original video is first submitted to a filtering processing (*e.g.*, Gaussian filter or Median filter). Then a subtraction is performed between the original and the filtered video, resulting in the noise residual video. Given that the highest responses representing the noise are concentrated on the abscissa and ordinate axes of the logarithm of the Fourier spectrum, visual rhythms [135, 136] are constructed to capture temporal information of the spectrum video sampling the central horizontal lines or central vertical lines of each frame and concatenating the sampling lines in a single image, called horizontal or a vertical visual rhythm. Then the grey-level co-occurrence matrices (GLCM) [137], LBP and HOG can be calculated on the visual rhythm as the texture features, and fed into a SVM or Partial Least Square (PLS). Furthermore, a more sophisticated method, based on Bag-of-Visual-Word model [138], similar to the VQ [124] used in [63], was also applied to extract the mid-level descriptor base on the low-level features, *e.g.*, LBP and HoG, extracted from the Fourier spectrum.

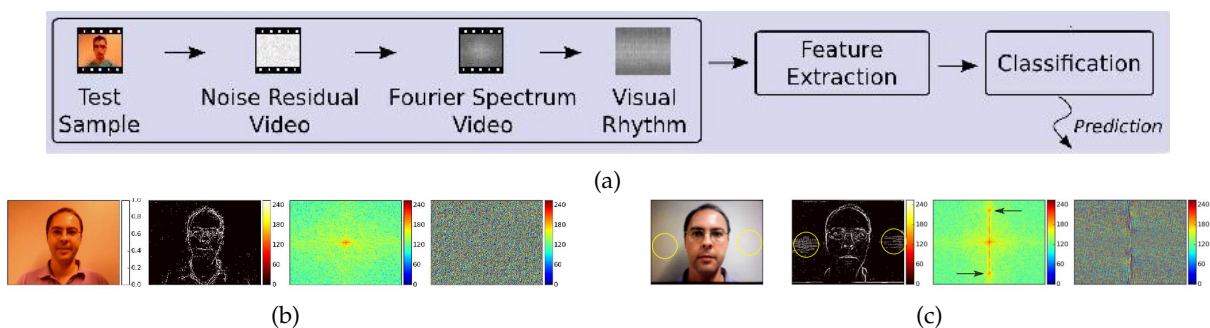


Figure 20: Illustration of the noise residual video and the visual rhythm based approach as shown in [78, 76] : (a) The framework of the visual rhythm based approach. (b) Example of valid access. (c) Example of a frame from a video replay attack. For (b) and (c), from left to right: original frame, residual noise frame, magnitude spectrum and phase spectrum. In (c), the yellow circles in the original image and its corresponding residual noise frame highlight the Moiré effect. The black arrows in the magnitude spectrum show the impact of the Moiré effect on the Fourier spectrum.

Also in 2012 and 2015 [76, 77], Tirunagari *et al.* [23] proposed to represent the dynamic characteristics of a video by a single image using Dynamic Mode Decomposition (DMD) [139]. Instead of sampling central lines in a video spectrum and concatenating them in a single frame (as in visual rhythms), the proposed approach selects the most representative frame in a video generated from the original video by applying DMD in the spatial space. DMD, similarly to Principal Component Analysis (PCA), is based on eigenvalues but, contrary to PCA, it can capture the motion in videos. The LBP feature of the DMD image is then calculated, and fed into a SVM for face PAD.

Xu *et al.* [79] first proposed to apply **deep learning** to learn the spatio-temporal features of a video for face PAD in 2015. More specifically, they proposed an architecture based on Long Short-Term Memory (LSTM) and CNN networks. As shown in Figure 21, several CNNs-based branches with only two convolutional layers are used. Each branch is used to extract the spatial texture features of one frame. These frames are sampled from the input video using a certain time step. Then, the LSTM units are connected at the end of each CNN branch to learn the temporal relations between frames. Finally, all the outputs of the LSTM units are connected to a softmax layer that gives the final classification of the input video for face PAD. Like several researchers before them, the authors also observed that using the scaled image of an original detected face including more background information can help the face PAD.

In 2019, Yang *et al.* [80] proposed a Spatio-Temporal Anti-Spoofing Network (STASN) to detect photo and video replay PAs. STASN consists of three modules: Temporal Anti-Spoofing Module (TASM), Region Attention Module (RAM), and Spatial Anti-Spoofing Module (SASM). The proposed TASM is composed of CNN and LSTM units, to learn the temporal features of the input video. One significant contribution is that, instead of using local regions with predefined locations as in [63, 52], STASN uses K local regions of the image selected automatically by RAM and TASM based on attention mechanism. These regions are then fed into SASM (*i.e.* a CNNs with K branches) for learning spatial texture features. STASN has

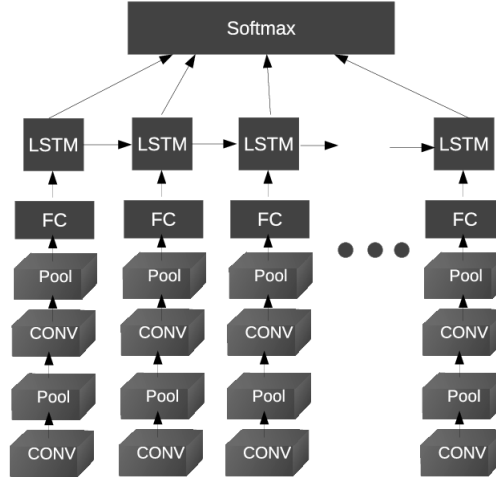


Figure 21: LSTM-CNN architecture used in [79] for face PAD.

significantly improved the performance for face PAD, especially in terms of the generalization capacity shown in cross-database evaluation scenarios (see Section 4).

2.4 3D geometric cue-based methods

3D geometric cue-based PAD methods use 3D geometric features to discriminate between a genuine face with a 3D structure that is characteristic of a face, and a 2D planar PA (e.g. a photo or video replay attack). The most widely used 3D geometric cues are the 3D shape reconstructed from the 2D image captured by the RGB camera, and the facial depth map, *i.e.* the distance between the camera and each pixel in the facial region. The two following sub-sections discuss the approaches based on these two cues, respectively.

2.4.1 3D shape-based methods

In 2013, Wang *et al.* [81] proposed a 3D shape-based method to detect photo attacks, in which the 3D facial structure is reconstructed from 2D facial landmarks [140] detected using different viewpoints [55, 141]. As shown in Figure 22, the reconstructed 3D structures of a real face and a planar photo are different. In particular, the reconstructed 3D structure from a real face profile preserves its 3D geometric structure. In contrast, the reconstructed structure of a planar photo in a profile view is only a line showing the photo's edge. The concatenation of the 3D coordinates of the reconstructed sparse structure are used as 3D geometric features, and fed into a SVM for face PAD. A drawback of this approach is that it requires multiple viewpoints, and cannot be used from a single image; using not enough key frames can lead to inaccuracies in the 3D structure reconstruction. Moreover, it is susceptible to inaccuracies in the detection of facial landmarks.

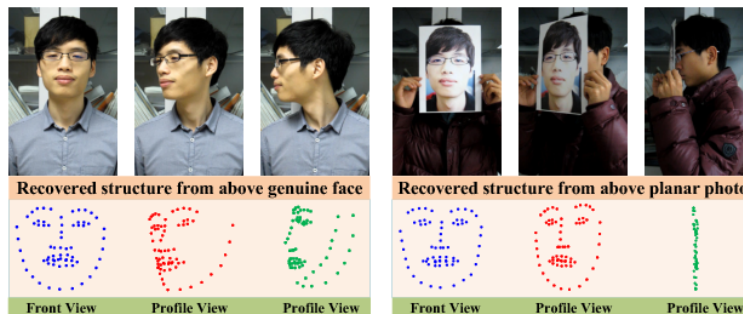


Figure 22: Illustration of reconstructed sparse 3D structures of genuine face (left) and photo attack (right) [81].

2.4.2 Pseudo-depth map-based methods

The depth map is defined as the distance of the face to the camera. Obviously, when using specific 3D sensors, the depth map can be captured directly. However, in this survey we focus on approaches that can be applied using GCDs, that do not usually embed 3D sensors. But, thanks to the significant progress in the computer vision area, especially with the **deep learning** technology, it is possible to get a good reconstruction / estimation of a depth map from a single RGB image [142, 143, 144]. Such reconstructions are called *pseudo-depth maps*. Based on the pseudo-depth map of a given image, the different PAD methods can be designed to discriminate between genuine faces and planar PAs.

In 2017, Atoum *et al.* [38] first proposed a depth-map based PAD method to detect planar face PADs, *e.g.*, printed photo attack and video replay attacks. The idea is to use the fact that the depth map of an actual face has varying height values in the depth map, whereas planar attacks' depth maps are constant (see Figure 23), to distinguish between real 3D faces and planar PAs. In this work, an 11-layer fully connected CNN [145] whose parameters are independent of the size of input face images is proposed, to estimate the depth map of a given image. The ground truth of depth maps was estimated using a state-of-the-art 3D face model fitting algorithm [143, 146, 147] for real faces, while it was set to zero for planar PAs, as shown in Figure 23. Finally, the estimated depth maps are fed to a SVM (pre-trained using the ground truth) to detect planar face PAs.

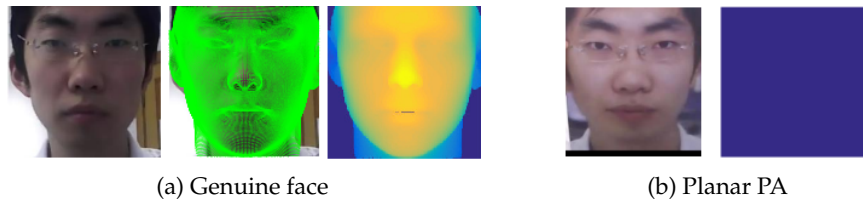


Figure 23: (a) A real face image, the fitted 3D face model and the depth map of the real 3D face; (b) A planar PA and its ground-truth depth map. The yellow/blue colors in the depth map represent respectively a closer/further point to the camera. Figure extracted from [38].

In 2018, Wang *et al.* [82] extended the single frame-based depth-map PAD method in [38] to videos, by proposing Face Anti-Spoofing Temporal-Depth networks (FAS-TD). FAS-TD networks are used to capture the motion and depth information of a given video. By integrating Optical Flow guided Feature Block (OFFB) and Convolution Gated Recurrent Units (ConvGRU) modules to a depth-supervised neural network architecture, the proposed FAS-TD can well capture short-term and long-term motion patterns of real faces and planar PAs in videos. The proposed FAS-TD further improved the performance of the depth-map based PAD methods using a single frame as [38, 34] and achieved state-of-the-art performances.

Since the pseudo-depth map approaches are very effective for detecting planar PAs, pseudo-depth maps are often used in conjunction with other cues, in multiple cues-based PAD method. Also, as pseudo-depth maps are among the most recently introduced cues for face PAD, they are extensively used in the most recent approaches. These two points are further detailed in the following Section 2.5, and in Section 2.6, respectively.

2.5 Multiple cue-based methods

Multi-modal systems are intrinsically more difficult to spoof than uni-modal systems [148, 62]. Some attempts to counterfeit face spoofing therefore combine methods based on different modalities, such as visible infrared [22], thermal infrared [149] or 3D [25] signals. However, the fact that such specific hardware is generally unavailable in most GCDs prevents these multi-modal solutions to be integrated into most existing face recognition systems. In this work, we focus on the multiple cues-based methods that use only images acquired using RGB cameras.

Such multiple cues-based methods combine liveness cues, texture cues and/or 3D geometric cues, to address the detection of various types of face PAs. In general, late fusion is used to merge the scores obtained from the different cues, to determine if the input image corresponds to a real face, or not.

2.5.1 Fusion of liveness cue and texture cues

In this section, the motion cue is used, as a liveness cue, in conjunction with different texture cues.

In 2017, Pan *et al.* [84] proposed to jointly use eye-blinking detection and the texture-based scene context matching for face PAD. The Conditional Random Field (CRF)-based eye-blinking model proposed by the same authors [43] (see page 9) is used to detect eye blinking. Then, a texture-based method is proposed, to check the coherence between the background region and the actual background (reference image). The reference image is acquired by taking a picture of the background, without the user being present. If the attempt is a real face presentation, the background region around the face in the reference image and the input image should theoretically be identical. Contrarily, if a video or re-captured photo (printed or displayed on a screen) is presented before the camera, then the background region around the face should be different between the reference image and the input image. To perform comparison between the input image's background and the reference image, LBP features are extracted from several fiducial points selected using the DoG function [150], and used to calculate the χ^2 distance as the scene matching score. If an imposter is detected by either the motion cue or the texture cue, then system will refuse access. This method has some limitations for real-life applications, as the camera should be fixed, and the background should not be monochrome.

The combination of motion cues and texture cues was widely used both in the first and second competitions on counter measures to 2D face spoofing attacks [151, 152] (held respectively in 2011 and in 2013). In the first competition, three of six teams used multiple cues-based methods. The AMILAB team used jointly face motion detection, face texture analysis and eye-blinking detection in their solution (and the sum of weighted classification scores obtained by SVMs). The CASIA team also considered three different cues: motion cue, noise-based texture cue and face-background dependency cue. The UNICAMP team combined eye blinking, background-dependency and micro-texture of an image sequence. In the second competition, the two teams that obtained the best performances (CASIA and LNMIT) used multiple-cues based methods. CASIA proposed an approach based on the early (feature-based) fusion of motion and texture cues, whereas LNMIT combined LBP, 2D FFT and face background consistency features [153] into a single feature vector, used as the input of Hidden Markov support vector machines (HM-SVMs) [154].

In 2016, Feng *et al.* [85] first integrated image quality measures (see page 16) as a static texture cue and motion cues in a neural network. Three different cues: Shearlet-based image quality features (SBIQF) [155, 156], a facial motion cue based on dense optical flow [157] and a scenic motion cue are manually extracted and fed into the neural network (see Figure 24). The neural network had been pre-trained, and was then fine-tuned on the existing for face PAD datasets.

2.5.2 Fusion of liveness and 3D geometric cues

In 2018, Liu *et al.* [34] proposed to use a two-stream CNN-RNN for fusing the remote PhotoPlethysmography (rPPG) cue and the pseudo-depth map cue for face PAD (see Figure 25). This approach uses the fully-connected CNN proposed in [38] to estimate the depth map (see page 21). Besides, a bypass connection is used to fuse the features from different layers, as in ResNet [4]. This work was the first one to proposed using RNN with LSTM units to learn the rPPG signal features based on the feature maps learned using CNNs. The estimation of the depth maps was calculated in advance using CNNs, whereas the depth maps' ground truth was estimated in advance using [143, 146, 147], and the rPPG ground-truth was generated as described in Section 2.2.2. The authors also designed a non-registration layer to align the input face to a frontal face, as the input of the RNN, for estimating the rPPG signal features. Instead of designing a binary classifier, the face PAs are then detected by thresholding a score computed based on the weighted quadratic sum of the estimated depth map of the last frame of the video, and the estimated rPPG signal features.

2.5.3 Fusion of texture and 3D geometric cues

In 2017, Atoum *et al.* [38] proposed to integrate patch-based texture cues and pseudo depth-map cues in a two-stream CNNs for face PAD. The pseudo-depth map estimation, that aims at extracting the holistic features of an image, has been described in Section 2.4.2. The patch-based CNNs stream (with 7 layers) focused on the image's local features. The local patches, with fixed size, are randomly extracted from the input image. The label of a patch extracted from a real face is set to 1, whereas the label of the patch of

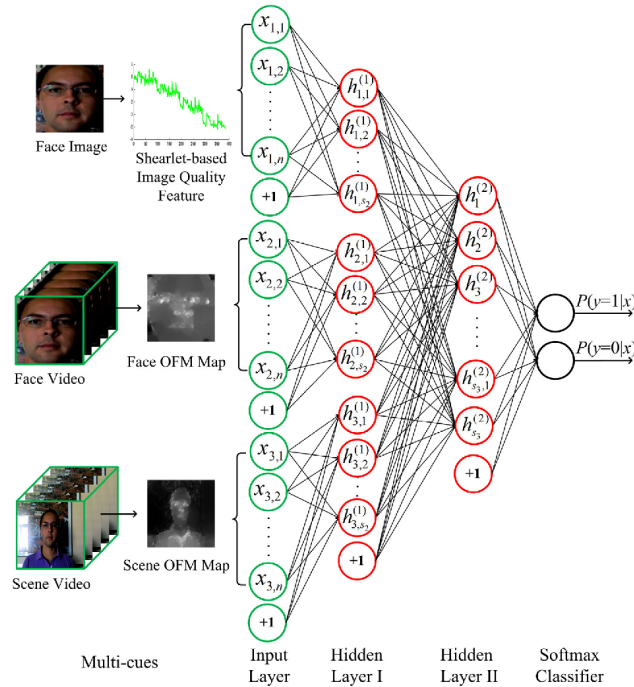


Figure 24: The flowchart of the multiple cues-based face PAD method using neural networks as shown in [85].

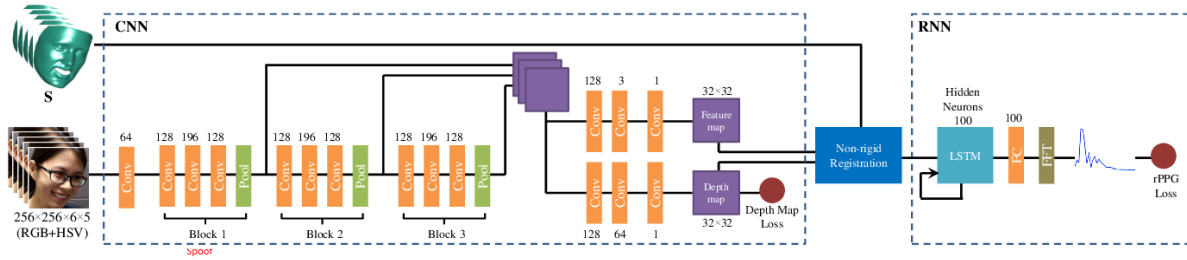


Figure 25: The proposed rPPG and depth-based CNN-RNN architecture for face PAD as shown in [34].

extracted from a PA is set 0. Then, the randomly extracted patches with their labels are used to train the patch-based CNNs stream with the softmax loss. Using the patch-level input cannot only increase the number of training samples, but also force the CNNs to explore the spoof-specific local discriminative information spreading in the entire face region. Finally, the two streams' scores are weighted to sum up as the final score to determine if the input image is a real face, or a PA. As in [33, 37, 61], the authors also proposed to jointly use the HSV/YCbCr image with the RGB image, as the input of the networks.

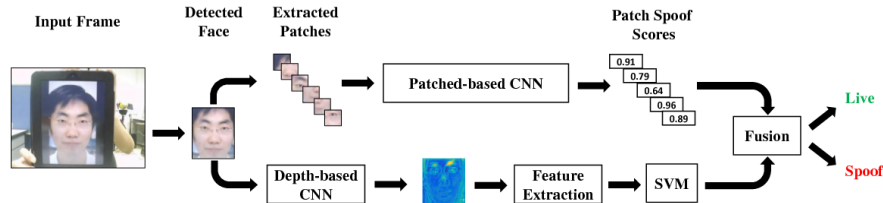


Figure 26: The architecture of the proposed patch-based and depth-based CNNs for face PAD as shown in [38].

2.6 New trends in PAD methods

In this section, we describe the methods that constitute the leading edge of face PAD methods based on RGB cameras. Thanks to the development of deep learning, especially in the computer vision domain, not only the face anti-spoofing detection performance has been significantly boosted, but also many new ideas have been introduced. These new ideas either relied on:

- the proposal of new cues to detect the face artifact (*e.g.* the pseudo-depth maps described in section 2.4.2);
- learning the most appropriate neural networks architectures for face PAD (*e.g.* using Neural Architecture Search (NAS) (see hereafter Section 2.6.1));
- address the generalization issues, especially towards types of attacks that are not (or insufficiently) represented in the learning dataset. Generalization issues can be (at least partially) addressed using zero / few shot learning (see Section 2.6.2) and/or domain adaptation and adversarial learning (see section 2.6.3).

The remainder of this section aims at presenting the two latter new trends more in details.

2.6.1 Neural Architecture Search (NAS) based PAD methods

In the last few years, deep neural networks have gained great success in many areas, such as speech recognition [158], image recognition [2, 159] and machine translation [160, 161]. The high performance of deep neural networks is heavily dependent on the adequation between their architecture and the problem at hand. For instance, the success of models like Inception [162], ResNets [4], and DenseNets [132], demonstrate the benefits of intricate design patterns. However, even with expert knowledge, determining which design elements to weave together generally requires extensive experimental studies [163]. Since the neural networks are still hard to design *a priori*, Neural Architecture Search (NAS) has been proposed to design the neural networks automatically based on reinforcement learning [164, 165], evolution algorithm [166, 167] or gradient-based methods [168, 169]. Recently, NAS has been applied to several challenging computer vision tasks, such as face recognition [170], action recognition [171], person re-identification [172], object detection [173] and segmentation [174]. However, NAS has just started being applied to face PAD.

In 2020, Yu *et al.* [83] first proposed to use NAS to design a neural network for estimating the depth map of a given RGB image for face PAD. The gradient-based DARTS [168] and Pc-DARTS [169] search methods were adopted to search the architecture of cells forming the network backbone for face PAD. Three levels of cells (low-level, mid-level, and high-level) from the three blocks of CNNs in [34] (see Figure 25) were used for the search space. Each block has four layers, including three convolutional layers and one max-pooling layer, and is represented as a Directed Acyclic Graph (DAG), with each layer as a node.

The DAG is used to present all the possible connections and operations between the layers in a block, as shown in Figure 27. Instead of directly using the original convolutional layers as in [34], the authors proposed to use Central Difference Convolutional (CDC) layers, in which the sample values in local receptive field regions are subtracted to the value of the central position, similarly to LBP. Then the convolution operation is based on the local receptive field region with gradient values.

A Multiscale Attention Fusion Module (MAFM) is also proposed for fusing low, mid and high levels CDC features via spatial attention [83]. Finally, the searched optimal architecture of the networks for estimating the depth map of a given image is shown in Figure 28. Rather than [34] fusing multiple cues in the CNNs, this work only estimated the depth map of an input image to employ face PAD by thresholding the mean value of the predicted depth map.

2.6.2 Zero/few-shot learning based PAD methods

Thanks to the significant development of deep learning, most state-of-the-art face PAD methods show promising performance in intra-database tests on the existing public datasets [38, 34, 83] (see Section 4). Nevertheless, the generalization to cross-dataset scenarios is still challenging, in particular due to the possible presence in the test set of face PAs that were not represented (or under-represented) in the training dataset [175, 66, 33]. One possible solution is to collect a large-scale dataset to include as much diverse spoofing attack types as possible. But, as detailed below in Section 3, unlike other problems such as face

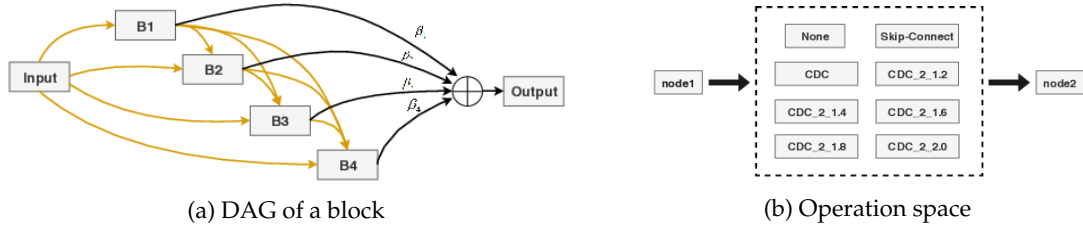


Figure 27: The search space of NAS for forming the network backbone for face PAD as shown in [83]. (a) Directed Acyclic Graph (DAG) of a block. Each node represents a layer in the block, and the edge is the possible information flow between layers. (b) Operation space, listing the possible operations between layers (8 operations were defined in [83]).

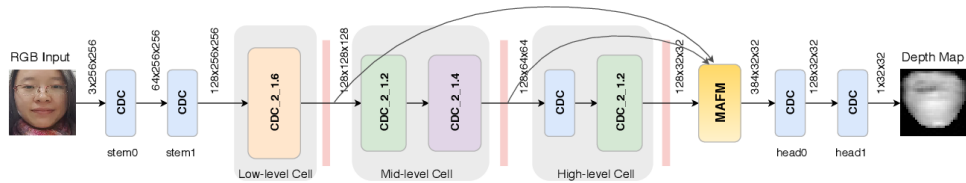


Figure 28: The architecture of the NAS-based backbone for depth map estimation as shown in [83].

recognition where it is relatively easy to collect massive public dataset from the Internet, the images / videos of spoofing artifacts re-captured by a biometric system are quite rarely available on Internet. Therefore, several research teams are currently investigating another solution, that consists in leveraging zero / few-shot learning to detect the previously unseen face PAs. This problem has been named Zero-Shot Face Anti-spoofing (ZSFA) in [88].

In 2017, Arashloo *et al.* [87] first addressed unseen attack detection, as an anomaly detection problem where real faces constitute the positive class, and used to train a **one-class classifier** such as one-class SVM [176].

In the same spirit, in 2018, Nikisins *et al.* [86] also used one-class classification. But, they used one-class Gaussian Mixed Models (GMM) to model the distribution of the real faces, in order to detect unseen attacks. Also, contrary to [87], they trained their model not only using one dataset, but aggregating three publicly available datasets (*i.e.* Replay-Attack [40], Replay-Mobile [40] and MSU MFSD [33], *c.f.* section 3).

The abovementioned methods used only samples of genuine faces to train one-class classifiers, whereas in practice, known spoof attacks might also provide valuable information to detect previously unknown attacks.

This is why, in 2019, Liu *et al.* [88] proposed a CNNs-based Deep Tree Network (DTN), in which 13 attack types covering both impersonation and obfuscation attacks are analyzed. First, they clustered the known PAs into eight semantic sub-groups using unsupervised tree learning, and they used them as the eight leaf nodes of the DTN (see Figure 29). Then, Tree Routing Unit (TRU) is learned to route the known PAs to the appropriate tree leaf (*i.e.* sub-group) based on the features of known PAs learned by the tree nodes (*i.e.* Convolutional Residual Unit (CRU)). In each leaf node, a Supervised Feature Learning (SFL) module, consisting of a binary classifier and a mask estimator, is employed to discriminate between spoofing attacks. The mask estimation is similar to the depth map estimation as in the same authors' previous work [34] (see page 11). Unseen attacks can then be discriminated based on the estimated mask, and the score of a binary softmax classifier.

2.6.3 Domain adaption based PAD methods

As detailed above, improving the generalization ability of existing face PAD methods is one of the greatest challenges nowadays. To mitigate this problem, Pereira *et al.* [74] first proposed to combine multiple databases to train the model (see page 19). This is the most intuitive attempt towards improving generalization of the earned models. However, even by combining all the existing datasets, it is impossible

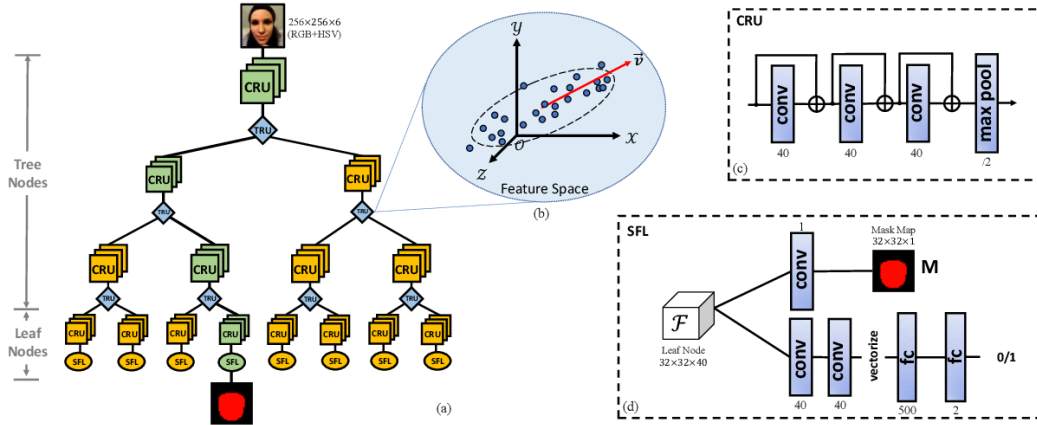


Figure 29: The proposed Deep Tree Network (DTN) architecture as shown in [88]. (a) Overall structure of the DTN. A tree node consists of a Convolutional Residual Unit (CRU) and a Tree Routing Unit (TRU), whereas a leaf node consists of a CRU and a Supervised Feature Learning (SFL) module. (b) Tree Routing Unit (TRU) assigns the feature learned by CRU to a given child node, based on an eigenvalue analysis similar to PCA. (c) Structure of each Convolutional Residual Unit (CRU). (d) Structure of the Supervised Feature Learning (SFL) in the leaf nodes.

to collect attacks from all possible domains (*i.e.* with every possible device, and in all possible capture environments), to train the model. However, even though printed photo or video replay attacks from unseen domains may differ greatly from the source domain, they all are based on paper or video screen as PAI [89]. Thus, if there exists a generalized feature space underlying the observed multiple source domains and the (hidden but related) target domain, then domain adaptation can be applied [177, 178].

In 2019, Shao *et al.* [89] first applied a domain adaption method based on adversarial learning [179, 180] to tackle face PAD. Under an adversarial learning schema, N discriminators were trained to help the feature generator produce generalized features for each of the N specific domains, as shown in Figure 30. Triplet loss is also used, to enhance the learned generalized features to be even more discriminative, both within a database (intra-domain) and among different databases (inter-domain). To apply face PAD, the learned generalized features are also trained to estimate the depth map of a given image as in [34], and to classify the image based on a binary classifier trained by softmax loss. This approach shows its superiority, when increasing the number of source domains, for learning generalized features. Indeed, in contrast, the previous methods without domain adaption such as LBP-TOP [72] or [34] cannot effectively improve the model's generalization capacity, even when using multiple source datasets for training.

3 Existing face anti-spoofing datasets and their major limitations

3.1 Some useful definitions

Face PAD (anti-spoofing) datasets consist of two different kinds of documents (files), in the form of photos or videos:

- The set of "**genuine faces**", that contains photos or videos of the genuine users' faces (authentic faces of the alive genuine users).
- The set of "**PA documents**", containing photos or videos of the PAI (printed photo, video replay, 3D mask, etc.)

Figure 31 illustrates the data collection procedure for constructing face anti-spoofing databases.

In face anti-spoofing datasets, genuine faces and PAIs (presented by imposters) are generally captured using the same device. This device is playing the role of the biometric system's camera in real-life authentication

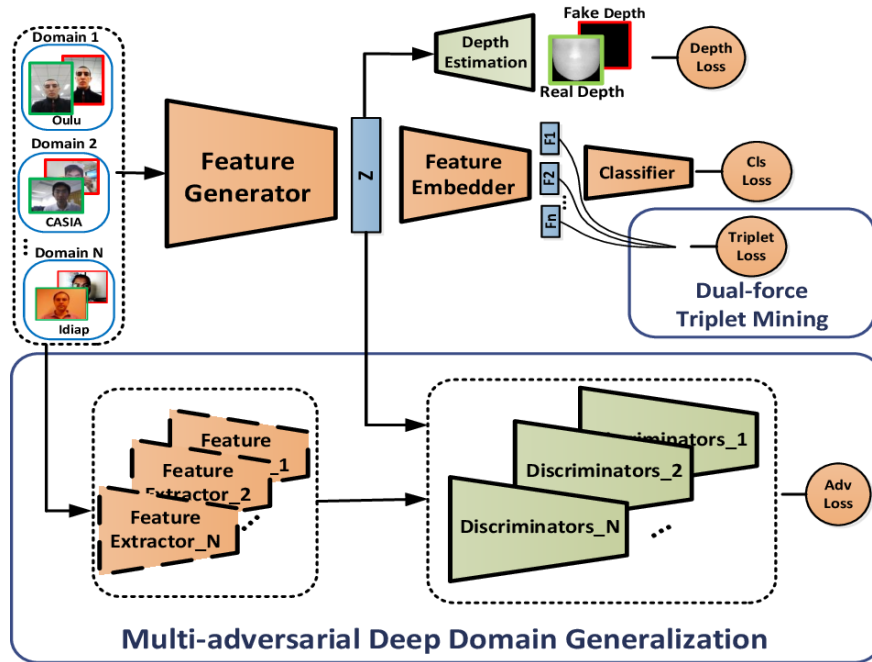


Figure 30: Overview of the domain adaption approach based on adversarial learning for face anti-spoofing, as shown in [89]. Each discriminator is trained to help the generalized features (learned by the generator from multiple source domains) to better generalize on their corresponding source domain. The depth loss, triplet loss and the classification loss ("Cls loss") are then used to enhance the ability to discriminate any kind of PAs.

applications; we therefore chose to call it "**biometric system acquisition device**". For genuine faces and 3D mask attacks, only the biometric system acquisition device is used to capture the data.

But, for printed photo and video replay attacks, another device is used to create the PAI (photo or video) from a genuine face's data. We call this device "**PA acquisition device**". It has to be noted that the PA acquisition device is in general different from the biometric system acquisition device. Some authors use the term of "**re-capture**", as the original data is first collected using the PA acquisition device, then presented on a PAI, and then re-captured using the biometric system acquisition device.

It has to be noted that, for photo display attacks and video replay attacks, the PAI itself can also be yet another electronic device. But, in general, only its screen is used, for displaying the PAI to the biometric system acquisition device. Of course, there could be datasets where the PAI also plays the role of PA acquisition device, but this is not the case in general. Indeed, it is not the case in most real-life applications, where the imposter generally does not have control over the PA acquisition device (*e.g.* photos or videos found on the web).

For printed photo attacks, paper-crafted mask attacks, and 3D mask attacks, yet another device is used: a printer. For photo attacks as well as paper-crafted mask attacks, usual (2D) printers are used. The printer's characteristics, as well as the quality of the paper used, can greatly affect the quality of the PAI, and therefore the chances of success of the attack. For 3D mask attacks, a 3D printer is used; its characteristics, as well as the material used (*e.g.* silicone or hard resine), and its thickness, also have an impact on the attack's chances of success.

The devices used for each dataset's collection are detailed in Table 2 and Section 3.4, together with a detailed description of these datasets. But, before that, in the remainder of this section, we successively give a brief overview of the existing datasets (Section 3.2) and describe their main limitations (Section 3.3).

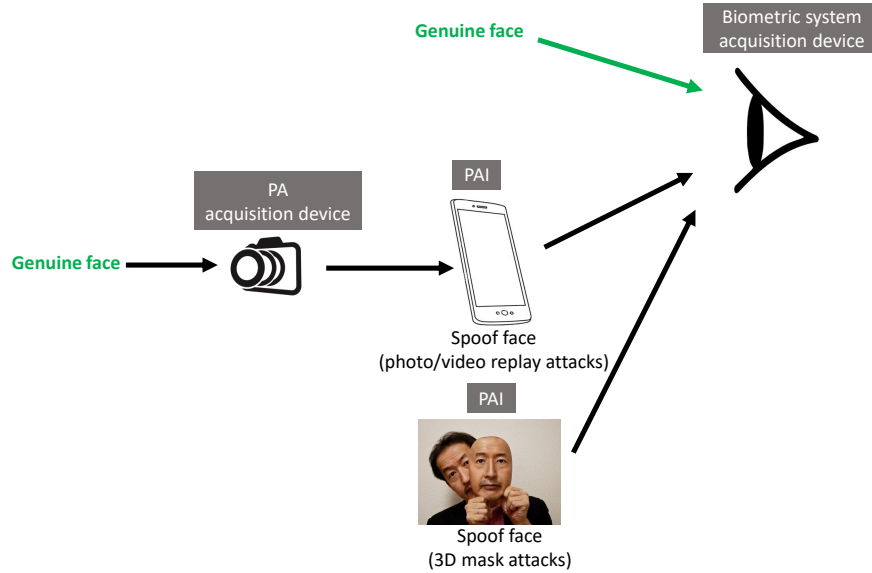


Figure 31: Diagram illustrating the data collection procedure for constructing face anti-spoofing databases.

3.2 Brief overview of the existing datasets

The early studies of face PAD, such as [45, 46, 47, 43, 84, 57], are mostly based on private datasets. Such private datasets being quite limited, both in volume and diversity of attack types, makes it very difficult to fairly compare the different approaches.

The first public dataset was proposed in 2008 by Tan *et al.* [53]. The dataset is named NUAA and it contains examples of photo attacks. The NUAA dataset enabled the researchers to compare the results of their methods on the same benchmark. Later on, respectively in 2011 and 2012, Anjos *et al.* publicly shared the datasets PRINT-ATTACK [148] (containing photo attacks) and its extended version REPLAY-ATTACK [39] (containing video replay attacks as well). Quickly, these two datasets were widely adopted by the research community, and so was one of the most challenging datasets: CASIA-FASD [19], that has also been published in 2012, and contains photo/video replay attacks, but with more diversity in the PAs, PAIs and video resolutions. Later on, several other similar datasets have been shared publicly for photo and video replay attacks, such as MSU-MFSD [33], MSU-USSA [67], OULU-NPU [181], SiW [34] and the very recent multi-modal dataset CASIA-SURF [182]. These datasets contain more diverse spoofing scenarios, such as MSU-MFSD [33] which first introduced a mobile phone scenario, MSU-USSA [67] which used celebrities' photos from the Internet to increase the mass of data, OULU-NPU [181] which focused on the attacks using mobile phones, and SiW [34] which contains faces with various poses, illumination and facial expressions.

The first public 3D mask attack dataset is 3DMAD [59], which includes texture maps, depth maps and point clouds together with the original images. Note that the depth maps and point clouds were collected by the Kinect 3D sensor rather than generic RGB cameras.

All the above-mentioned datasets have been created under controlled environments, *i.e.* mostly indoors and with controlled illumination conditions, face poses *etc.*

Although UAD [77] has collected videos from both indoors and outdoors, and with a relatively large number of subjects, the dataset is no longer publicly available. Although MSU-USSA [67] contains a set of genuine faces captured under more diverse environments (including celebrities' images collected from the Internet by [183]), the PAs always took place in controlled indoors conditions. Even the latest CASIA-SURF dataset [182], which is so far the largest multi-modal face anti-spoofing dataset with 1000 subjects, contains only images collected in the same well-controlled conditions.

Therefore, public datasets are still far from reproducing real-world applications in a realistic way. This is probably due to the difficulty of collecting impostors' PAs and PAIs in the wild. As a consequence, examples of PAs are generally acquired manually, which is a time-consuming and draining work. And, creating a

large-scale dataset for face anti-spoofing in the wild, covering realistically various real-world applicative scenarios, is still a challenge. To circumvent these challenges, some researchers use data augmentation techniques to create synthetic (yet realistic) images of PAs [80].

A summarized overview of the existing public face anti-spoofing datasets using only generic RGB cameras is provided in Table 2. More precisely, for each dataset, Table 2 gives (in columns) its release's year (*Year*), the number of subjects it contains (*# Subj*), the ethnicity of the subjects in the dataset (*Ethnicity*), the type of PA represented in the dataset (*PA type(s)*), the number and type(s) of documents provided in the dataset as the cumulated number of genuine attempts and PAs (*Document # & type(s)*), the PAI(s) used (*PAI*), the head pose(s) in the set of genuine faces (*Pose*), whether or not there are facial expression variations in the genuine faces dataset (*Expressions*), the biometric system acquisition device(s) for capturing both the genuine attempts and, in case of an attack, the PAI (*Biometric system acquisition device*) and the PA acquisition device that is possibly used to create the PAI (*PA acquisition device*).

3.3 Major limitations of the existing datasets

Given the acquisition difficulties mentioned above, the existing face PAD datasets are (compared to other face related problems) still limited not only in terms of volume, but also in terms of diversity regarding the types of PAs, PAIs and acquisition devices used for genuine faces, PAs, and possibly PAIs. In particular, as of today, there is still no public large-scale face PAD in the wild, whereas there are several such datasets for face recognition.

This hinders the development of effective face PAD methods. It partly explains why, compared to other face-related problems, such as face recognition, the performances of the current face PAD methods are still below the requirements of most real-world applications (especially in terms of their generalization ability).

Of course, this is not the only reason: as detailed earlier in this paper, face PAD is a very challenging problem. But, because all data-driven (learning-based) methods' performances – including hand-crafted feature-based methods and more recent deep learning-based methods – are largely affected by the learning dataset's volume and diversity [130, 131, 184, 185], the lack of diversity in the datasets contribute to the limited performances of the current face PAD methods.

More details about these datasets, including discussion about their advantages and drawbacks, are provided in the remainder of this section.

Table 2: A summary of public face anti-spoofing datasets based on generic RGB camera.

Database	Year	# Subj.	Ethnicity	PA type(s)	Document # & type(s) Images (I) / Videos (V)	PAI	Pose	Expression	Biometric system acquisition device	PA acquisition device
NUAA [53]	2010	15	Asian	Printed photos Warped photos	5105/7509 (I)	A4 paper	Frontal	No	Webcam(640x480)	Webcam(640x480)
PRINT-ATTACK [148]	2011	50	Caucasian	Printed photos	200/200 (V)	A4 paper	Frontal	No	Macbook Webcam (320x340)	Cannon PowerShot SX150 (12.1 MP)
CASIA-FASD[19]	2012	50	Asian	Printed photos Warped photos Cut photos Video replay	200/450 (V)	Copper paper iPad 1 (1024x768)	Frontal	No	Sony NEX-5 (1280x720) USB Camera (640x480)	Sony NEX-5 (1280x720) Webcam (640x480)
REPLAY-ATTACK [39]	2012	50	Caucasian 76% Asian 22% African 2%	Printed photos Photo display 2x video replays ^a	200/1000 (V)	A4 paper iPad 1 (1024x768) iPhone 3GS (480x320)	Frontal	No	Macbook Webcam (320x340)	Canon PowerShot SX 150 (12.1MP) iPhone 3GS
3DMAD [59, 25]	2013	17	Caucasian	2x 3D masks ^b	170/85 (V)	Paper-crafted mask Hard resin mask (ThatsMyFace.com)	Frontal	No	Kinect (RGB camera) (Depth sensor)	—
MSU-MFSD [33]	2015	35	Caucasian 70% Asian 28% African 2%	Printed photos 2x video replays	110/330 (V)	A3 paper iPad Air (2048x1536) iPhone 5s (1136x640)	Frontal	No	Nexus 5 (built-in camera software 720x480) Macbook Air (640x480)	Cannon 550D (1920x1088) iPhone 5s (1920x1080)
MSU-RAFS [60]	2015	55	Caucasian 44% Asian 53% African 3%	Video replays	55/110 (V)	Macbook (1280x800)	Frontal	No	Nexus 5 (rear: 3264x2448) iPhone 6 (rear: 1920x1080)	The biometric system acquisition devices used in MSU-MFSD, CASIA-FASD, REPLAY-ATTACK.
UAD [77]	2015	404	Caucasian 44% Asian 53% African 3%	7x video replays	808/16,268 (V)	7 display devices	Frontal	No	6 different cameras (no mobile phone) (1366x768)	6 different cameras (no mobile phone) (1366x768)
MSU-USSA [67]	2016	1,140	Diverse set (from web faces database from the [183])	Printed photos Photo display 3x video replays	1,140/9,120 (V)	White paper (11x8.5 paper) Macbook (2080x1800) Nexus 5 (1920x1080) Tablet (1920x1200)	Frontal	Yes	Nexus 5 (front: 1280x960) (rear: 3264x2448) iPhone 6 (rear: 1920x1080)	Same as MSU-RAFS Cameras used to capture celebrities' photos are unknown.
OULU-NPU [181]	2017	55	Caucasian 5% Asian 95%	Printed photos 2x video replays	1,980/3,960 (V)	A3 glossy paper Dell display (1280x1024) Macbook (2560x1600)	Frontal	No	Samsung Galaxy S6 (rear: 16 MP)	Samsung Galaxy S6 (front: 5 MP) HTC Desire EYE (front: 13 MP) MEIZU X5 (front: 5 MP) ASUS Zenfone Selfi (front: 13 MP) Sony XPERIA C5 (front: 13 MP) OPPO N3 (front: 16 MP)
SiW [34]	2018	165	Caucasian 35% Asian 35% African American 7% Indian 23%	Printed photos (high/low-quality photos) 4x video replays	1,320/3,300 (V)	Printed paper (High/low-quality) Samsung Galaxy S8 iPhone 7 iPad Pro PC screen(Asus MB168B)	[−90°, 90°]	Yes	Camera (1920x1080)	Camera (1920x1080) Camera (5184x3456)
CASIA-SURF [182]	2019	1000	Asian	Flat-cut/Warped-cut photos (eyes, nose, mouth)	3,000/18,000 (V)	A4 paper	[−30°, 30°]	No	RealSense (RGB camera) (1280x720) (Depth sensor) (640x480) (IR sensor) (640x480)	RealSense (RGB camera) (1280x720) (Depth sensor) (640x480) (IR sensor) (640x480)

^a $x \times$ **video replays** denotes x types of video replay attacks with different PAIs.

^b $2 \times$ **3D masks** denotes two types of 3D masks attacks: paper-crafted masks and hard resin mask.

3.4 Detailed description of the existing datasets

In this section, we provide a detailed description of all the datasets mentioned in Table 2.

NUAA Database [53] is the first publicly available face PAD dataset for printed photo attacks. It includes some variability in the PAs, as the photos are moved/distorted in front of the PA acquisition device as follows:

- 4 kinds of translations: vertical, horizontal, toward the sensor, toward the background
- 2 kinds of rotations: along the horizontal axis and along the vertical axis (in-depth rotation)
- 2 kinds of bending: along the horizontal and vertical axis (inward and outward)

A generic webcam is used for recording the genuine face images. Fifteen subjects are enrolled in the database, and each subject is asked to avoid eye blinking and to keep a frontal pose, with neutral facial expression. The attacks are performed by using printed photographs (either on photographic paper or A4 paper printed by a usual color HP printer). The dataset is divided into two separate subsets, for training and test. The training set contains 1743 genuine face images and 1748 PAs impersonating 9 genuine users. The test set contains 3362 genuine samples and 5761 PAs. Viola-Jones detector [186] is used to detect the faces in the images, and the detected faces are aligned/normalized according to the eyes locations detected by [187]. The face images are then resized to 64×64 pixels. Extracts from the NUAA database are shown in Figure 32.



Figure 32: NUAA (from left to right): five different photo attacks.

PRINT-ATTACK Database [148] is the second proposed public dataset, including photo-attacks impersonating 50 different genuine users. The data is collected in two different conditions: controlled and adverse. In controlled conditions, the scene background is uniform and the light of a fluorescent lamp illuminates the scene, while in adverse conditions, the scene background is non-uniform and day-light illuminates the scene. A MacBook is used to record video clips of the genuine faces and the PAs. To capture the photos used for the attack, a 12.1 megapixel *Canon PowerShot SX150 IS* camera is used. These photos are then printed on plain A4 paper using a *Triumph-Adler DCC 2520* color laser printer. Video clips of about 10 seconds are captured for each PA, under two different scenarios: hand-based attacks and fixed-support attacks. In hand-based attacks, the impostor holds the printed photos using their own hands, whereas in fixed-support attacks, the impostors stick the printed photos to the wall so they do not move/shake during the PA. Finally, 200 genuine attempts and 200 PA video clips are recorded. The 400 video clips are then divided into three subsets: training, validation, and testing. Genuine identities (real identities or impersonated identities) in each subset were chosen randomly but with no overlap. Extracts of the PRINT-ATTACK dataset are shown in Figure 33.

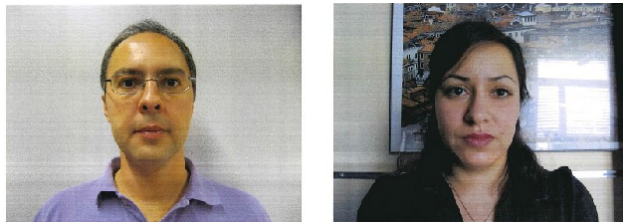


Figure 33: PRINT-ATTACK (from left to right): photo attack under controlled and adverse scenarios.

CASIA-FASD Database [19] is the first publicly available face PAD dataset that provides both printed photo and video replay attacks. The CASIA-FASD database is a spoofing attack database which consists of three types of attacks: warped printed photos (which simulates paper mask attacks), printed photos with cut eyes and video attacks (motion cue such as eye blinking is also included). Each real face video and spoofing attack video is collected in three different qualities: low, normal and high quality. The high-quality video has a high resolution 1280×720 , and the low/normal quality video has the same resolution 640×480 . However, the low and normal quality is defined empirically by the perceptual feeling rather than strict quantitative measures. The whole database is split into a training set (containing 20 subjects) and a testing set (containing 30 subjects). Seven test scenarios are designed considering three different image qualities, three different attacks (warped/cut photo attack and video replay attack) and the overall test combining all the data. Examples of CASIA-FASD database are shown in Figure 34.



Figure 34: CASIA-FASD (from left to right): real face, two warped/cut photo attacks and a video replay attack.

REPLAY-ATTACK Database [39] is an addendum of the above-mentioned PRINT-ATTACK database [148] proposed by the same team. Compared to the PRINT-ATTACK database, REPLAY-ATTACK adds two more attacks, which are Phone-Attack and Tablet-Attack. The Phone-Attack uses an iPhone screen to display the video or photo attack, and the Tablet-Attack uses an iPad screen to display high resolution (1024×768) digital photos or videos. Thus, REPLAY-ATTACK database can be used to evaluate photo attacks using printed photo or screens, and video replay attacks. The number of video clips for spoof attacks is increased from 200 to 1000, for 50 identities (subjects). The dataset is divided into training, validation and test sets. REPLAY-ATTACK database also offers an extra subset as the enrollment videos for 50 genuine clients, to be used for evaluating the vulnerabilities of a face recognition system without face PAD is vulnerable towards various various types of attacks. Examples of REPLAY-ATTACK database are shown in Figure 35.



Figure 35: REPLAY-ATTACK (from left to right): real face, video replay attack, photo displayed on screen, printed photo attack.

3DMAD Database [59, 25] is the first public face anti-spoofing database for 3D mask attack. Previous databases contain attacks performed with 2D artifacts (*i.e.* photo or video) that are in general unable to fool face PAD systems relying on 3D cues. In this database, the attackers wear customized 3D facial masks made out of a hard resin (manufactured by ThatsMyFace.com) of a valid user to impersonate the real access. It is worth mentioning that paper-craft mask files are also provided in this dataset. The dataset contains a total of 255 videos of 17 subjects. For each access attempt, a video was captured using the *Microsoft Kinect for Xbox 360*, which provides RGB data and depth information of size 640×480 at 30 frames per second. This dataset allows to evaluate both 2D and 3D PAD techniques, and also their fusion. It is divided into three sessions: two real access sessions recorded with a time delay and one attack session captured by a single operator (attacker). Examples of 3DMAD database are shown in Figure 36.



Figure 36: 3DMAD (from left to right): paper-craft mask and 17 hard resin masks.

MSU-MFSD Database [33] is the first publicly available database to use mobile phones to capture real accesses. This database includes real access and attack videos for 55 subjects (among which 35 subjects in the public version: 15 subjects in the training set and 20 subjects in test set). The genuine faces were captured using two devices: a Google Nexus 5 phone using its front camera (720×480 pixels) and a MacBook Air using its built-in camera (640×480 pixels). The Canon 550D SLR camera (1920×1088) and iPhone 5S (rear camera 1920×1080) are used to capture high-resolution pictures or videos (for photo attacks and video replay attacks). The printed high-resolution photo is played back using an iPhone 5S as PAI, and high definition (HD) (1920×1088) video-replays (captured on a Canon 550D SLR) are played back using an iPad Air. Examples of MSU-MFSD database are shown in Figure 37.



Figure 37: MSU-MFSD (from left to right): genuine face, video replay attacks respectively displayed on iPad and iPhone, and printed photo attack.

MSU-RAFS Database [60] is an extension of MSU-MFSD [33], CASIA-FASD [19] and REPLAY-ATTACK [39], where the video replay attacks are generated by replaying (on a MacBook) the genuine face videos in MSU-MFSD, CASIA-FASD and REPLAY-ATTACK. Fifty-five videos are genuine face videos from MSU-MFSD (captured by using the front camera of a Google Nexus 5), while 110 (2×55) videos are videos replay attacks, captured using the built-in rear camera of a Google Nexus 5 and the built-in rear camera of an iPhone 6, and re-played using a MacBook as PAI. In addition, 100 genuine face videos from CASIA-FASD and REPLAY-ATTACK were both used as genuine face videos, and used to generate 200 video replay attacks by replaying these genuine face videos using a MacBook as PAI. During the attack, the average standoff of the smartphone camera (used by the biometric system) from the screen of the MacBook was 15 cm, which assured that replay videos do not contain the bezels (edges) of the MacBook screen. Unlike the previously described databases, MSU-RAFS is constructed using existing genuine face videos (without having control over the biometric system acquisition devices used). Therefore, in this dataset, the biometric system acquisition devices used for capturing the genuine face videos generally differ from the devices used for capturing the PAs. Thus, there is a risk to introduce a bias when evaluating methods based on this dataset only.

Examples of MSU-RAFS database are shown in Figure 38.

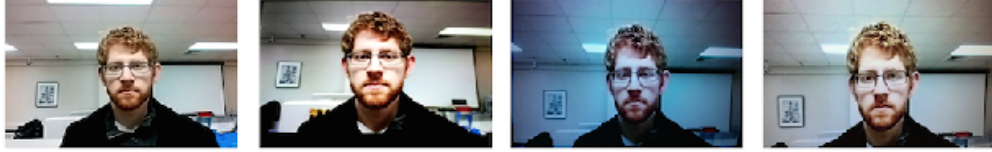


Figure 38: MSU-RAFS (from left to right): genuine face, PA from MSU-MFSD, (attacks using a MacBook (as a PAI) to replay the genuine attempts from MSU-MFSD, captured by different devices (respectively Google Nexus 5 and iPhone 6).)

UAD Database [77] is the first database to collect the data both indoors and outdoors. It is also much bigger than the previous databases, both in terms of the number of subjects (440 subjects) and the number of videos (808 for training/16,268 for testing). All videos have been recorded at full-HD resolution, but subsequently cropped and resized to 1366×768 pixels. The dataset includes real access videos collected using six different cameras. For each subject, two videos are provided, both using the same camera but under different ambient conditions. Spoof attack videos corresponding to a given subject have also been captured using the same camera as for his / her real access videos. The video replay attacks have been displayed using seven different electronic monitors. But, this database seems to be no longer publicly available nowadays. Examples of MSU-RAFS database are shown in Figure 39.



Figure 39: UAD (from left to right): video replay attacks, captured outdoors (first and second images) and indoors (last three images).

MSU-USSA Database [67] can be regarded as an extension of the MSU-RAFS [60], proposed by the same authors. There are two sub-sets in the database: 1) following the same idea as for MSU-RAFS, the first subset consists of 140 subjects from REPLAY-ATTACK [39] (50 subjects), CASIA-FASD [19] (50 subjects) and MSU-MFSD [33] (40 subjects); 2) the second subset consists of 1000 subjects taken from the web faces database collected in [183], containing images of celebrities taken under a variety of backgrounds, illumination conditions and resolutions. Only a single frontal face image of each celebrity is retained. Thus, MSU-USSA database contains color face images of 1140 subjects, where the average resolution of genuine face images is 705×865 . Two cameras (front and rear cameras of a Google Nexus 5 smartphone) have been used to collect 2D attacks using four different PAIs (laptop, tablet, smartphone and printed photos), resulting in a total of 1140 genuine faces and 9120 PAs. Just like MSU-RAFS, MSU-USSA has not captured the genuine face videos with the same device used for capturing the PAs. Thus, there is a risk to introduce a bias when evaluating methods based on this dataset only. Examples of MSU-USSA database are shown in Figure 40.



Figure 40: MSU-USSA (from left to right): spoof faces re-captured from the celebrity dataset [183].

OULU-NPU Database [181] is a more recent dataset (introduced in 2017), that contains PAD attacks acquired with mobile devices. In most previous datasets, the images were acquired in constrained conditions. On the other hand, this database contains a variety of motion, blur, illumination conditions, backgrounds and head

poses. The database includes data corresponding to 55 subjects. The front cameras of 6 different mobile devices have been used to capture the images included in this dataset. The images have been collected under three separate conditions (environment / face artifacts / acquisition devices), each corresponding to a different combination of illumination and background. Presentation attacks include printed photo attacks created using two printers, as well as video replay attacks using two different display devices. Four protocols are proposed for methods benchmarking (see Section 4.3 for more details). In total, the dataset is composed of 4950 real accesses and attack videos. Examples of OULU-NPU database are shown in Figure 41.



Figure 41: Extracts of the OULU-NPU dataset.

SiW Database [34] is the first database to include faces spoofing attacks with both various labelled poses and facial expressions. This database consists of 1320 genuine access videos captured from 165 subjects and 3300 attack videos. Compared to the above mentioned databases, it includes subjects from a wider variety of ethnicities, *i.e.* Caucasian (35%), Indian (23%), African American (7%) and Asian (35%). Two kinds of print (photo) attacks and four kinds of video replay attacks have been included in this dataset. Video replay attacks have been created using four spoof mediums (PAIs): two smartphones, a tablet, and a laptop. Four different sessions corresponding to different head poses / camera distances, facial expressions and illumination conditions were collected, and three protocols were proposed for benchmarking (see Section 4.3 for more details). Examples of the SiW database are shown in Figure 42.



Figure 42: SiW: genuine access (top) and PA (bottom) videos with different poses, facial expressions and illumination conditions (for genuines accesses), and PAI devices (for PAs).

CASIA-SURF Database [60] is currently the largest face anti-spoofing dataset containing multi-modal images, *i.e.* RGB (1280×720), Depth (640×480) and Infrared (IR) (640×480) images, of 1000 subjects in 21,000 videos. Each sample includes one live (genuine) video clip and six spoof (PA) video clips under different types of attacks. Six different photo attacks are included in this database: flat/warped printed photos where different regions are cut from the printed face. During the dataset capture, the genuine users and imposters were required to turn left or right, move up or down, walk towards or away from the camera (imposters holding the printed color photo on an A4 paper). The face angle was only limited to 300 degrees. Imposters stood within a range of 0.3 to 1.0 meter from the camera. The RealSense SR300 camera was used to capture the RGB, Depth and Infrared (IR) images. The database is divided into three subsets for training, validation and testing. In total, there are 300 subjects and 6300 videos (2100 for each modality) in the training set, 100 subjects and 2100 videos (700 for each modality) in the validation set, and 600 subjects and 12600 videos (4200 for each modality) in the testing set. Examples of the CASIA-SURF database are shown in Figure 43.

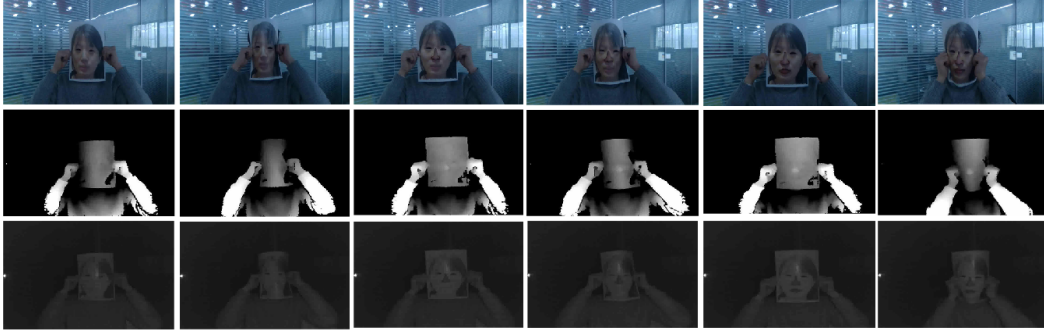


Figure 43: Extract from CASIA-SURF showing six photo attacks in each 3 modalities: RGB (top), depth (middle) and IR (middle).

4 Experimental evaluation

In this section, we present a comprehensive evaluation of the approaches for face PAD detailed in Section 2. By doing so, our objective is to investigate the strengths and weaknesses of the different types of methods, in order to draw future research directions for face PAD, so as to make face authentication less vulnerable to imposters. We first present (in Section 4.1) the evaluation protocol, then (in Section 4.2), the evaluation metrics, and finally (in Section 4.3) the experimental results.

4.1 Evaluation protocol

In this section, we present the protocol we used to compare experimentally the different face PAD methods. In the early studies of face anti-spoofing detection, there was no uniform protocol to train and evaluate the face PAD methods. In 2011, a first standard protocol was proposed by Anjos *et al.* [148] in order to fairly compare the different methods. In 2017, a second standard protocol was proposed by Boulkenaf *et al.* [181] based on an ISO/IEC standard. On top of the evaluation metrics to be used, these protocols address mainly two aspects: 1) how to divide the database; 2) what kinds of tests should be conducted for the evaluation, *e.g.* intra-database and inter (cross)-database tests.

a) Dataset division The protocol proposed by Anjos *et al.* [148] is widely used when the evaluation is based on PRINT-ATTACK [148], REPLAY-ATTACK [39], OULU-NPU [181], 3DMAD [59, 25] and / or CASIA-SURF [182]. This protocol relies on the division of the dataset into three subsets: training, validation set and test sets (respectively for training, tuning the model’s parameters and assessing the performances of the tuned model).

Other databases, such as CASIA-FASD[19], MSU-MFSD [33] and SiW [34], only consist of two independent subsets: training and test subsets. In this case, either a small part of the training set is used as a validation set, or cross-validation is used, to tune the model’s parameters. In some datasets (such as OULU-NPU and SiW), the existence of different sessions explicitly containing different capture conditions allowed the authors to propose refined protocols for evaluation (as detailed above in Section 3.4), and can also be used for dataset division.

b) Intra-database vs inter-database evaluation An intra-database evaluation protocol uses only a single database to both train and evaluate the PAD methods. But, as the current databases are still limited in terms of variability, intra-database evaluations can be subject to overfitting, and therefore report biased (optimistic) results.

The inter-database test, proposed by Pereira *et al.* [74] for evaluating the generalization abilities of a model, consists in training the model on a certain database, and then evaluate it on a separate database. Although inter-database tests (or cross-database tests) aim to evaluate the model’s generalization abilities, it is important to note that the evaluation performances are still affected by the distribution of the two datasets. Indeed, if the two datasets’ distribution is close (*e.g.* if the same PAIs or spoof acquisition devices were used),

the inter-database test will also report optimistic results. For instance, for a given model, inter-database evaluations between PRINT-ATTACK and MSU-MFSD (using the same MacBook camera to acquire the images) result in much better performances than inter-database evaluations between CASIA-FASD and MSU-MFSD (where CASIA-FASD uses a USB camera with very different from a MacBook), as reported in [33].

4.2 Evaluation metric

To compare different PAD methods, Anjos *et al.* [148] proposed in 2011 to use Half Total Error Rate (HTER) as an evaluation metric. HTER combines the False Rejection Rate (FRR) and the False Acceptance Rate (FAR) as follows:

$$HTER = \frac{FRR + FAR}{2} \quad (1)$$

where the FAR and FRR are respectively defined as:

$$FAR = \frac{FP}{FP + TN} \quad (2)$$

$$FRR = \frac{FN}{FN + TP} \quad (3)$$

with TP, FP, TN and FN respectively corresponding to the numbers of True Positives, False Positives, True Negatives and False Negatives. TP, FP, TN and FN are calculated using the model parameters achieving Equal Error Rate (EER) on the validation set (parameters for which $FRR=FAR$). It can be noted that EER is also often used for assessing the model's performance on the validation and training subsets.

But, since 2017 and the work proposed in [181], the performance is most often reported using the metrics defined in the standardized ISO/IEC 30107-3 metrics [16]: Attack Presentation Classification Error Rate (APCER) and Bona Fide Classification Error Rate (BPCER) (also called Normal Presentation Classification Error Rate (NPCER) in some research papers). These two metrics correspond respectively to the False Acceptance Rate (FAR) and the False Rejection Rate (FRR), but for obtaining APCER, the FAR is computed separately for each PAI / type of attack, and APCER is defined as the highest FAR (*i.e.* the FAR of the most successful type of attack). Similar to HTER, the Average Classification Error Rate (ACER) is then defined as the mean of APCER and BPCER using the model parameters achieving EER on the validation set:

$$ACER = \frac{APCER + NPCER}{2} \quad (4)$$

On top of the HTER and ACER scalar values, the Receiver Operating Characteristic (ROC) curve and the Area Under the Curve (AUC) are also commonly used to evaluate the PAD method's performance. The two latter have the advantage that they can provide a global evaluation of the model's performances over different values of the parameter set.

4.3 Experimental results

In this part, we compare some of the face PAD methods detailed in Section 2 on the public benchmarks presented in Section 3, following the intra-database and inter/cross-database protocols described above. Among the more than 50 methods presented in Section 2, we selected here the most influential methods and / or the ones that are among the most characteristic of their type of approach (following the typology presented in Section 2.1 and Figure 3, page 6 and used for the methods presentation in Section 2). To compare the performances of the different methods, we used the metrics EER, HETR, APCER, BPCER and ACER described above. The results we report are extracted from the original papers introducing the methods, *i.e.* we did not re-develop all these methods to perform the evaluation ourselves. As a consequence, some values might be missing (and are then noted as "-") in the following tables. Another point that is important to note is that we chose to focus our analysis on the type of features used. Of course, depending on the method, the type of classifier used (or the neural network architecture, for end-to-end deep learning methods) might also have an impact on the overall performance. But, we consider that, for each method, the authors have chosen to use the most effective classifier / architecture, and that therefore the overall performances they report are largely representative of the descriptive and discriminative capabilities of the features they use.

4.3.1 Intra-database evaluation on public benchmarks

Table 3 and Table 4 respectively show the results of the intra-database evaluation on the CASIA-FASD and REPLAY-ATTACK datasets. Compared to static texture feature-based methods such as DoG, LBP-based methods, dynamic texture-based methods such as LBP-TOP, Spectral Cubes [78] and DMD [23] are more effective on both benchmarks. However, the static features learned using CNNs can boost the performance significantly, and sometimes even outperform dynamic texture hand-crafted features. For instance, even the earliest CNNs-based method with static texture feature [36] has shown a superior performance in terms of HTER than almost all previously introduced state-of-the-art methods based on hand-crafted features. It was the first time that the deep CNNs showed its potential for face PAD. Later, researchers proposed more and more models learning static or dynamic features based on deep CNNs such as LSTM-CNN [79], DPCNN [69] and Patch-based CNN [38], that has achieved the state-of-the-art performances on both the CASIA-FASD and REPLAY-ATTACK datasets. Besides, we can see both in Table 3 and Table 4 that Patch-Depth CNN [38], fusing different cues (*i.e.* texture cue and 3D geometric cue (depth map)), has shown its superiority over single cue-based methods using the same CNN, such as Patch cue-based CNN [38] or Depth cue-based CNN [38]. Indeed, their HTERs are respectively 2.27%, 2.85% and 2.52% on CASIA-FASD, and 0.72%, 0.86% and 0.75% on REPLAY-ATTACK. These results show the effectiveness of multiple cues-based methods that, by leveraging different cues, are able to effectively detect a wider variety of PA types.

As explained earlier in section 3.4, on OULU-NPU and SiW, the different protocols corresponding to different applicative scenarios were proposed.

More specifically, for the benchmark OULU-NPU, four protocols were proposed in [181]:

- Protocol 1 aims at testing the PAD methods under different environmental conditions (illumination and background);
- Protocol 2's objective is to test the generalization abilities of the methods learnt using different PAIs
- Protocol 3 aims at testing the generalization across the different acquisition devices (*i.e.* using Leave One Camera Out (LOCO) protocol to test the method over six smartphones);
- Protocol 4 is the most challenging scenario, as it combines the three previous protocols to simulate real-world operational conditions.

For SiW, three different protocols were proposed in [34]:

- Protocol 1 deals with variations in face pose and expression;
- Protocol 2 tests the model over different spoof mediums (PAIs) for video replay;
- Protocol 3 tests the methods over different PAs, *e.g.* learning from photo attacks and testing on video attacks, and *vice versa*.

From Table 5 and Table 6, one can see that, both for OULU-NPU and SiW and for all the evaluation protocols, the best methods are the 3D geometric cue methods using depth estimation. Furthermore, the architectures obtained using NAS with depth maps (*e.g.* CDCN++ [83]) has achieved state-of-the-art performances both on OULU-NPU and SiW .

Moreover, we can see that the protocols used on OULU-NPU for testing generalization abilities (protocols 2 and 3) are especially challenging. When considering the protocol defined to evaluate the performances in near real-world applicative conditions (protocol 4), the model's performance can degrade up to 25 times compared to "easier" protocols (*e.g.*, CDCN++'s ACER arises from 0.2% for protocol 1 to 5.0% for protocol 4). Similar results can be observed on SiW.

It indicates that the generalization across the scenarios is still a challenge for face PAD methods, even within the same dataset.

4.3.2 Cross-database evaluation on public benchmarks

Compared to the promising results shown in the intra-database test, the inter/cross-database test results are still way worse than most real-world applications requirements. Several databases have been adopted to perform cross (inter)-database evaluation, such as CASIA-FASD *vs.* MSU-MFSD [33], MSU-USSA *vs.*

Table 3: Evaluation of various face PAD methods on CASIA-FASD.

Method	Year	Feature	Cues	EER(%)	HTER(%)
DoG [19]	2012	DoG	Texture (static)	17.00	-
LBP [39]	2012	LBP	Texture (static)	-	18.21
LBP-TOP [73]	2014	LBP	Texture (dynamic)	10.00	-
Yang <i>et al.</i> [36]	2014	CNN	Texture (static)	4.92	4.95
Spectral Cubes [78]	2015	FourrierSpectrum +codebook	Texture (dynamic)	14.00	-
DMD [23]	2015	LBP	Texture (dynamic)	21.80	-
Color texture [37]	2015	LBP	Texture (HSV/static)	6.20	-
Moire pattern [60]	2015	LBP+SIFT	Texture (static)	-	0
LSTM-CNN [79]	2015	CNN	Texture (dynamic)	5.17	5.93
Color LBP [61]	2016	LBP	Texture (HSV/static)	3.20	-
Fine-tuned VGG-Face [69]	2016	CNN	Texture (static)	5.20	-
DPCNN [69]	2016	CNN	Texture (static)	4.50	-
Patch-based CNN [38]	2017	CNN	Texture (static)	4.44	3.78
Depth-based CNN [38]	2017	CNN	Depth	2.85	2.52
Patch-Depth CNN [38]	2017	CNN	Texture+Depth	2.67	2.27

Table 4: Evaluation of various face PAD methods on REPLAY-ATTACK.

Method	Year	Feature	Cues	EER(%)	HTER(%)
LBP [39]	2012	LBP	Texture (static)	13.90	13.87
Motion Mag [75]	2013	HOOOF	Texture (dynamic)	-	1.25
LBP-TOP [73]	2014	LBP	Texture (dynamic)	7.88	7.60
Yang <i>et al.</i> [36]	2014	CNN	Texture (static)	2.54	2.14
Spectral Cubes [78]	2015	FourrierSpectrum +codebook	Texture (dynamic)	-	2.80
DMD [23]	2015	LBP	Texture (dynamic)	5.30	3.80
Color texture [37]	2015	LBP	Texture (HSV/static)	0.40	2.90
Moire pattern [60]	2015	LBP+SIFT	Texture (static)	-	3.30
Color LBP [61]	2016	LBP	Texture (HSV/static)	0.10	2.20
Fine-tuned VGG-Face [69]	2016	CNN	Texture (static)	8.40	4.30
DPCNN [69]	2016	CNN	Texture (static)	2.90	6.10
Patch-based CNN [38]	2017	CNN	Texture (static)	2.50	1.25
Depth-based CNN [38]	2017	CNN	Depth	0.86	0.75
Patch-Depth CNN [38]	2017	CNN	Texture+Depth	0.79	0.72

REPLAY-ATTACK / CASIA-FASD / MSU-MFSD [67]. However, most researchers have reported their cross-database evaluation results using REPLAY-ATTACK *vs.* CASIA-FASD [74, 34, 83, 80, 70, 82], since the important differences between these two databases introduce a great challenge for cross-database testing.

Table 7 reports the results of cross-database tests between REPLAY-ATTACK and CASIA-FASD. Although the use of deep learning methods significantly improves the generalization between different datasets, there is still a large gap compared to the intra-database results. Especially if we train the model on REPLAY-ATTACK and then test the trained model on CASIA-FASD, even the best methods can only achieve at best a 29.8% HTER.

Moreover, all the PAD methods based on hand-crafted features show weak generalization abilities. For instance, the HTER of LBP-based methods based on RGB image (such as basic LBP [32] and LBP-TOP [74]) is about 60%. However, the LBP in HSV/YCbCr color space shows a comparable or even better generalization ability than some deep learning-based methods (*e.g.*, the method in [61] achieves respectively a 30.3% and 37.7% HTER when trained on CASIA-FASD and REPLAY-ATTACK). It is noteworthy that the multiple cues-based method Auxiliary [34], by fusing depth map and rPPG cues, achieves a good generalization even in the most difficult cross-database tests. For instance, when trained on REPLAY-ATTACK and tested on CASIA-FASD, it achieves slightly better HTER (28.4%) than the latest method CDCN++ [83] based on NAS. (29.8%, see Table 7). This demonstrates that the multiple cues-based methods, when using different cues that

Table 5: Evaluation of various face PAD methods on OULU-NPU.

Protocol	Method	Year	Feature	Cues	APCER(%)	BPCER(%)	ACER(%)
1	CPqD [175]	2017	Inception-v3 [188]	Texture (static)	2.9	10.8	6.9
1	GRADIANT [175]	2017	LBP	Texture (HSV/static)	1.3	12.5	6.9
1	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	1.6	1.6	1.6
1	FaceDs [70]	2018	CNN	Texture (Quality/static)	1.2	1.7	1.5
1	STASN [80]	2019	CNN+Attention	Texture (dynamic)	1.2	2.5	1.9
1	FAS_TD [82]	2019	CNN+LSTM	Depth	2.5	0.0	1.3
1	DeepPixBis [71]	2019	DenseNet [132]	Texture	0.8	0.0	0.4
1	CDCN [83]	2020	CNN	Depth	0.4	1.7	1.0
1	CDCN++ [83]	2020	NAS+Attention	Depth	0.4	0.0	0.2
2	MixedFASNet [175]	2017	DNN	Texture (HSV/static)	9.7	2.5	6.1
2	GRADIANT [175]	2017	LBP	Texture (HSV/static)	3.1	1.9	2.5
2	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	2.7	2.7	2.7
2	FaceDs [70]	2018	CNN	Texture (Quality/static)	4.2	4.4	4.3
2	STASN [80]	2019	CNN+Attention	Texture (dynamic)	4.2	0.3	2.2
2	FAS_TD [82]	2019	CNN+LSTM	Depth	1.7	2.0	1.9
2	DeepPixBis [71]	2019	DenseNet [132]	Texture (static)	11.4	0.6	6.0
2	CDCN [83]	2020	CNN	Depth	1.5	1.4	1.5
2	CDCN++ [83]	2020	NAS+Attention	Depth	1.8	0.8	1.3
3	MixedFASNet [175]	2017	DNN	Texture (HSV/static)	5.3 ± 6.7	5.30 ± 6.7	5.3 ± 6.7
3	GRADIANT [175]	2017	LBP	Texture (HSV/static)	2.6 ± 3.9	5.0 ± 5.3	3.8 ± 2.4
3	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	2.7 ± 1.3	3.1 ± 1.7	2.9 ± 1.5
3	FaceDs [70]	2018	CNN	Texture (Quality/static)	4.0 ± 1.8	3.8 ± 1.2	3.6 ± 1.6
3	STASN [80]	2019	CNN+Attention	Texture (dynamic)	4.7 ± 3.9	0.9 ± 1.2	2.8 ± 1.6
3	FAS_TD [82]	2019	CNN+LSTM	Depth	5.9 ± 1.9	5.9 ± 3.0	5.9 ± 1.0
3	DeepPixBis [71]	2019	DenseNet [132]	Texture	11.7 ± 19.6	10.6 ± 14.1	11.1 ± 9.4
3	CDCN [83]	2020	CNN	Depth	2.4 ± 1.3	2.2 ± 2.0	2.3 ± 1.4
3	CDCN++ [83]	2020	NAS+Attention	Depth	1.7 ± 1.5	2.0 ± 1.2	1.8 ± 0.7
4	Massy_HNU [175]	2017	LBP	Texture (HSV+YCbCr)	35.8 ± 35.3	8.3 ± 4.1	22.1 ± 17.6
4	GRADIANT [175]	2017	LBP	Texture (HSV/static)	5.0 ± 4.5	15.0 ± 7.1	10.0 ± 5.0
4	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	9.3 ± 5.6	10.4 ± 6.0	9.5 ± 6.0
4	FaceDs [70]	2018	CNN	Texture (Quality/static)	1.2 ± 6.3	6.1 ± 5.1	5.6 ± 5.7
4	STASN [80]	2019	CNN+Attention	Texture (dynamic)	6.7 ± 10.6	8.3 ± 8.4	7.5 ± 4.7
4	FAS_TD [82]	2019	CNN+LSTM	Depth	14.2 ± 8.7	4.2 ± 3.8	9.2 ± 3.4
4	DeepPixBis [71]	2019	DenseNet [132]	Texture (static)	36.7 ± 29.7	13.3 ± 14.1	25.0 ± 12.7
4	CDCN [83]	2020	CNN	Depth	4.6 ± 4.6	9.2 ± 8.0	6.9 ± 2.9
4	CDCN++ [83]	2020	NAS+Attention	Depth	4.2 ± 3.4	5.8 ± 4.9	5.0 ± 2.9

Table 6: Evaluation of various face PAD methods on SiW.

Protocol	Method	Year	Feature	Cues	APCER(%)	BPCER(%)	ACER(%)
1	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	3.58	3.58	3.58
1	STASN [80]	2019	CNN+Attention	Texture (dynamic)	-	-	1.0
1	FAS_TD [82]	2019	CNN+LSTM	Depth	0.96	0.50	0.73
1	CDCN [83]	2020	CNN	Depth	0.07	0.17	0.12
1	CDCN++ [83]	2020	NAS+Attention	Depth	0.07	0.17	0.12
2	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	0.57 ± 0.69	0.57 ± 0.69	0.57 ± 0.69
2	STASN [80]	2019	CNN+Attention	Texture (dynamic)	-	-	0.28 ± 0.05
2	FAS_TD [82]	2019	CNN+LSTM	Depth	0.08 ± 0.14	0.21 ± 0.14	0.14 ± 0.14
2	CDCN [83]	2020	CNN	Depth	0.00 ± 0.00	0.13 ± 0.09	0.06 ± 0.04
2	CDCN++ [83]	2020	NAS+Attention	Depth	0.00 ± 0.00	0.09 ± 0.10	0.04 ± 0.05
3	Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG	8.31 ± 3.81	8.31 ± 3.80	8.3 ± 3.81
3	STASN [80]	2019	CNN+Attention	Texture (dynamic)	-	-	12.10 ± 1.50
3	FAS_TD [82]	2019	CNN+LSTM	Depth	3.10 ± 0.81	3.09 ± 0.81	3.10 ± 0.81
3	CDCN [83]	2020	CNN	Depth	1.67 ± 0.11	1.76 ± 0.12	1.71 ± 0.11
3	CDCN++ [83]	2020	NAS+Attention	Depth	1.97 ± 0.33	1.77 ± 0.10	1.90 ± 0.15

are inherently complementary to each other, can achieve better generalization than face PAD models based

on single cues. But, from a very general perspective, improving the generalization abilities of the current face PAD methods is still a great challenge for face anti-spoofing.

Table 7: Cross-database testing between CASIA-FASD and REPLAY-ATTACK. The reported evaluation metric is HTER(%).

Method	Year	Feature	Cues	Train	Test	Train	Test
				CASIA-FASD	REPLAY-ATTACK	REPLAY-ATTACK	CASIA-FASD
LBP [32] ^a	2012	LBP	Texture (static)		55.9		57.6
Correlation 19 [189] ^a	2013	MLP	Motion		50.2		47.9
LBP-TOP [74]	2013	LBP	Texture (dynamic)		49.7		60.6
Motion Mag [75]	2013	HOOF	Texture+Motion		50.1		47.0
Yang <i>et al.</i> [36]	2014	CNN	Texture (static)		48.5		45.5
Spectral cubes [78]	2015	FourrierSpectrum+codebook	Texture (dynamic)		34.4		50.0
Color texture [37]	2015	LBP	Texture (HSV/static)		47.0		39.6
Color LBP [61]	2016	LBP	Texture (HSV/static)		30.3		37.7
Auxiliary [34]	2018	CNN+LSTM	Depth+rPPG		27.6		28.4
FaceDs [70]	2018	CNN	Texture (Quality/static)		28.5		41.1
STASN [80]	2019	CNN+Attention	Texture (dynamic)		31.5		30.9
FAS_TD [82]	2019	CNN+LSTM	Depth		17.5		24.0
CDCN [83]	2020	CNN	Depth		15.5		32.6
CDCN++ [83]	2020	NAS+Attention	Depth		6.5		29.8

^a Results taken from [74].

5 Discussion

From the evaluation results presented in the previous section 4.3, we can see that face PAD is still a very challenging problem. In particular, the performances of the current face PAD methods are still below the requirements of most real-world applications (especially in terms of generalization ability).

More precisely, the performances are acceptable when there is not too much variation between the conditions of the genuine faces capture for enrollment, and the genuine face / PA presentation for authentication (intra-database evaluation).

But, as:

- all hand-crafted features show a limited generalization ability, as they are not powerful enough to capture all the possible variations in the acquisition conditions;
- the features learned by deep / wide neural networks are of very high dimensions, compared to the limited size of the training data,

both types of features suffer from over-fitting, and therefore poor generalization capacity.

Therefore, learning features that are able to discriminate between a genuine face and any kind of PA, possibly under very different capture conditions, is still an open issue. This issue will be discussed in Section 5.1. Then, in Section 5.2, we discuss a less studied topic in the field of face PAD: how to detect obfuscation attacks.

5.1 Current trends and perspectives

As stated earlier, learning features that are distinctive enough to discriminate between genuine faces and various PAs, possibly in very different environments, is still an open issue. Of course, this kind of issues (related to the generalization abilities of data-driven models) are common in the field of computer vision, way beyond face PAD.

However, as collecting impostors' PAs and PAIs is nearly impossible in the wild, collecting / creating a face PAD dataset with sufficient samples and variability (not only regarding the capture conditions, but also regarding the different types of PA / PAIs) is still very time-consuming and costly (see section 3.2). It is indeed much easier to create a dataset for most object recognition tasks (*e.g.* face authentication).

In order to tackle all previously seen PAs, a current trend is to combine multiple cues (see section 2.5). But, due to the above-mentioned challenges in the dataset creation as well as the technological advances that ill-intentioned users can access to deploy increasingly sophisticated attacks, it might happen that the PAD method has to detect PAs that were not included in its training dataset. This problem, called "Unknown attack" previously, is especially challenging.

Beyond the current methods that try and use zero/few-shot learning approaches to tackle this problem, the question of learning features that are representative enough of "real" faces, so that they can discriminate between genuine faces and any kind of PAs, under any type of capture conditions, is still an open issue. This issue, for which some researchers have recently proposed solutions based on domain adaptation, will very probably raise a lot of attention from researchers in the coming years, especially with the emergence of ever more sophisticated DeepFakes.

5.2 Obfuscation face PAD

As stated in Section 1.2, two types of PAs are defined in the relevant ISO standard [16]: impersonation (spoofing) attacks, *i.e.* attempts of impostors to impersonate a genuine user, and obfuscation attacks, *i.e.* attempts for the impostor to hide her own identity.

Most current face PAD research focuses on the former type, *i.e.* impersonation spoofing, as it is the most frequent attack for biometric systems based on face recognition / authentication.

However, there are some applicative scenarios where obfuscation attacks are very important to detect. For instance, in law-enforcement applications based on video-surveillance, one of the main objectives is to be able to detect criminals, whereas the goals of criminals using obfuscation attacks is to remain unrecognized by the system.

As detailed in Figure 1 on page 3, obfuscation attacks may entail (possibly extreme) facial makeup, or occluding significant portions of the face using scarves, sunglasses, face masks, hats, etc. In some cases, the person deploying obfuscation attack may also use tricks that are usually used for impersonation attacks, *e.g.* by using a mask showing the face of a non-criminal.

To the best of our knowledge, so far the only dataset containing examples of obfuscation attacks is SiW-M. This dataset has been introduced in [88], where the authors have shown the effectiveness of extreme makeup for face obfuscation. One solution is to process the face image so as to synthetically "remove" the makeup, as in [190, 191].

More generally, given that, compared to impersonation attacks, obfuscation attacks are still less frequent, several research groups consider obfuscation attacks as a zero / few-shot PAD problem [88].

Even though obfuscation attack detection has been so far much less studied than impersonation attack detection, it is very likely that this topic will become more and more studied in the future, given the conjunction of several factors, such as the generalization of video-surveillance in public places, geo-political issues including risks of terrorist attacks in some regions of the world, and recent technological developments that allow researchers to tackle this problem.

6 Conclusion

In this survey paper, we have thoroughly investigated over 50 of the most influential face PAD methods that can work in scenarios where the user only has access to the RGB camera of Generic Consumer Devices. By structuring our paper according to a typology of face PAD methods based on the types of PA they are aiming to thwart, and in chronological order, we have shown the evolution in the face PAD area during the last two decades. This evolution covers a large variety of methods, from hand-crafted features to the most recent deep learning-based technologies such as Neural Architecture Search (NAS). Benefiting from the recent breakthroughs obtained by researchers in Computer Vision, thanks to the advent of deep learning, face PAD methods are getting ever more effective, and efficient. We have also gathered, summarized and detailed the most relevant information about a dozen of most widespread public datasets for face PAD.

Using these datasets as benchmarks, we have extensively compared different types of face PAD methods using common experimental protocol and evaluation metrics. This comparative evaluation allows us to point out which types of approaches are most effective, depending on the type of PA. More specifically, according to our investigation, texture-based methods which are also the most widely used PAD methods, and especially dynamic texture-based methods, are able to detect almost all types of PAs. Furthermore, the methods based on texture features learned using deep learning have significantly improved the state-of-the-art face PAD performances, compared to methods based on hand-crafted texture features. However, in general, high quality 3D mask attacks are still a great challenge for texture-based approaches. On the other hand, liveness-based methods or 3D geometric-based methods can achieve relatively better generalization capabilities, even though they are still vulnerable to video replay attacks, or complex illumination conditions. Multiple cues-based methods, by leveraging different cues for face PAD, are in general more effective for detecting various PAs. Nevertheless, the computational complexity of multiple-cues based methods is an issue needed to be considered for real-time applications. Partly because of the complexity of the face PAD problem, of the huge variability in the possible attacks and the lack of dataset that contains enough samples with sufficient variability, all current approaches are still limited in terms of generalization.

We have also identified some of the most prominent current trends in face PAD, such as combining approaches that aim at thwarting various kinds of attacks, or tackling previously unseen attacks. We have also provided some insights for future research and have listed the still open issues, such as learning features that are able to discriminate between genuine faces and all kinds of PAs.

References

- [1] Matthew Turk and Alex Pentland. Eigenfaces for recognition. *Journal of cognitive neuroscience*, 3(1):71–86, 1991.
- [2] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012.
- [3] Christian Szegedy, Wei Liu, and et al. Going deeper with convolutions. In *CVPR*, pages 1–9, 2015.
- [4] Kaiming He, Xiangyu Zhang, and et al. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [5] Yaniv Taigman, Ming Yang, and et al. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, pages 1701–1708, 2014.
- [6] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deeply learned face representations are sparse, selective, and robust. In *CVPR*, pages 2892–2900, 2015.
- [7] Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In *BMVC*, volume 1, page 6, 2015.
- [8] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015.
- [9] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphreface: Deep hypersphere embedding for face recognition. In *The CVPR*, volume 1, page 1, 2017.
- [10] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019.
- [11] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation from predicting 10,000 classes. In *CVPR*, pages 1891–1898, 2014.
- [12] Gary B Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007.
- [13] Lior Wolf, Tal Hassner, and Itay Maoz. Face recognition in unconstrained videos with matched background similarity. In *CVPR*, pages 529–534. IEEE, 2011.
- [14] Rodrigo de Luis-García, Carlos Alberola-López, Otman Aghzout, and Juan Ruiz-Alzola. Biometric identification systems. *Signal processing*, 83(12):2539–2557, 2003.
- [15] Luiz Souza, Luciano Oliveira, Mauricio Pamplona, and Joao Papa. How far did we get in face spoofing detection? *Engineering Applications of Artificial Intelligence*, 72:368–381, 2018.

- [16] Iso/iec jtc 1/sc 37 biometrics. information technology – biometric presentation attack detection – part 1: Frame-work. *International Organization for Standardization*, 2016.
- [17] Siqi Liu, Baoyao Yang, Pong C Yuen, and Guoying Zhao. A 3d mask face anti-spoofing database with real world variations. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 100–106, 2016.
- [18] Klaus Kollreider, Hartwig Fronthaler, and Josef Bigun. Verifying liveness by multiple experts in face biometrics. In *CVPR Workshops*, pages 1–6. Ieee, 2008.
- [19] Zhiwei Zhang, Junjie Yan, and et al. A face antispoofing database with diverse attacks. In *International Conference on Biometrics*, pages 26–31. IEEE, 2012.
- [20] <http://www.urmesurveillance.com/>.
- [21] Sébastien Marcel, Mark S Nixon, and Stan Z Li. *Handbook of biometric anti-spoofing*, volume 1. Springer, 2014.
- [22] Zhiwei Zhang, Dong Yi, Zhen Lei, and Stan Z Li. Face liveness detection by learning multispectral reflectance distributions. In *Face and Gesture 2011*, pages 436–441. IEEE, 2011.
- [23] Santosh Tirunagari, Norman Poh, David Windridge, Aamo Iorliam, Nik Suki, and Anthony TS Ho. Detection of face spoofing using visual dynamics. *IEEE transactions on information forensics and security*, 10(4):762–777, 2015.
- [24] Javier Galbally and Riccardo Satta. Three-dimensional and two-and-a-half-dimensional face recognition spoofing using three-dimensional printed models. *IET Biometrics*, 5(2):83–91, 2016.
- [25] Nesli Erdogmus and Sebastien Marcel. Spoofing face recognition with 3d masks. *IEEE transactions on information forensics and security*, 9(7):1084–1097, 2014.
- [26] Andrea Lagorio, Massimo Tistarelli, Marinella Cadoni, Clinton Fookes, and Sridha Sridharan. Liveness detection based on 3d face shape analysis. In *2013 International Workshop on Biometrics and Forensics (IWBF)*, pages 1–4. IEEE, 2013.
- [27] Stan Li, Dong Yi, Zhen Lei, and Shengcai Liao. The casia nir-vis 2.0 face database. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 348–353, 2013.
- [28] Sushil Bhattacharjee and Sébastien Marcel. What you can’t see can help you-extended-range imaging for 3d-mask presentation attack detection. In *2017 International Conference of the Biometrics Special Interest Group (BIOSIG)*, pages 1–7. IEEE, 2017.
- [29] Javier Hernandez-Ortega, Julian Fierrez, Aythami Morales, and Pedro Tome. Time analysis of pulse-based face anti-spoofing in visible and nir. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 544–552, 2018.
- [30] Dong Yi, Zhen Lei, Zhiwei Zhang, and Stan Z Li. Face anti-spoofing: Multi-spectral approach. In *Handbook of Biometric Anti-Spoofing*, pages 83–102. Springer, 2014.
- [31] Lin Sun, WaiBin Huang, and MingHui Wu. Tir/vis correlation for liveness detection in face recognition. In *International Conference on Computer Analysis of Images and Patterns*, pages 114–121. Springer, 2011.
- [32] Ivana Chingovska, Andre Anjos, and Sebastien Marcel. On the effectiveness of local binary patterns in face anti-spoofing. In *2012 BIOSIG - Proceedings of the International Conference of Biometrics Special Interest Group (BIOSIG)*, pages 1–7, 2012.
- [33] Di Wen, Hu Han, and Anil K Jain. Face spoof detection with image distortion analysis. *IEEE Transactions on Information Forensics and Security*, 10(4):746–761, 2015.
- [34] Yaojie Liu, Amin Jourabloo, and X. Liu. Learning deep models for face anti-spoofing: Binary or auxiliary supervision. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 389–398, 2018.
- [35] Tiago de Freitas Pereira, Andre Anjos, Jose Mario De Martino, and Sebastien Marcel. Can face anti-spoofing countermeasures work in a real world scenario? In *2013 International Conference on Biometrics (ICB)*, pages 1–8, 2013.
- [36] Jianwei Yang, Zhen Lei, and Stan Z Li. Learn convolutional neural network for face anti-spoofing. *arXiv preprint arXiv:1408.5601*, 2014.
- [37] Zinelabidine Boulkenafet, Jukka Komulainen, and Abdenour Hadid. Face anti-spoofing based on color texture analysis. In *2015 IEEE international conference on image processing (ICIP)*, pages 2636–2640. IEEE, 2015.

- [38] Yousef Atoum, Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. Face anti-spoofing using patch and depth-based cnns. In *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pages 319–328, 2017.
- [39] Ivana Chingovska, André Anjos, and Sébastien Marcel. On the effectiveness of local binary patterns in face anti-spoofing. In *International Conference of the Biometrics Special Interest Group (BIOSIG)*, pages 1–7. IEEE, 2012.
- [40] Artur Costa-Pazo, Sushil Bhattacharjee, Esteban Vazquez-Fernandez, and Sebastien Marcel. The replay-mobile face presentation-attack database. In *2016 International Conference of the Biometrics Special Interest Group (BIOSIG)*, pages 1–7. IEEE, 2016.
- [41] Raghavendra Ramachandra and Christoph Busch. Presentation attack detection methods for face recognition systems: A comprehensive survey. *ACM Computing Surveys (CSUR)*, 50(1):1–37, 2017.
- [42] Javier Hernandez-Ortega, Julian Fierrez, Aythami Morales, and Javier Galbally. Introduction to face presentation attack detection. In *Handbook of Biometric Anti-Spoofing*, pages 187–206. Springer, 2019.
- [43] Gang Pan, Lin Sun, and et al. Eyeblick-based anti-spoofing in face recognition from a generic webcam. In *ICCV*, pages 1–8. IEEE, 2007.
- [44] Lin Sun, Gang Pan, Zhaohui Wu, and Shihong Lao. Blinking-based live face detection using conditional random fields. In *International Conference on Biometrics*, pages 252–260. Springer, 2007.
- [45] Jiangwei Li, Yunhong Wang, Tieniu Tan, and Anil K Jain. Live face detection based on the analysis of fourier spectra. In *Biometric technology for human identification*, volume 5404, pages 296–303. International Society for Optics and Photonics, 2004.
- [46] Klaus Kollreider, Hartwig Fronthaler, and Josef Bigun. Evaluating liveness by face images and the structure tensor. In *Fourth IEEE Workshop on Automatic Identification Advanced Technologies (AutoID’05)*, pages 75–80. IEEE, 2005.
- [47] Klaus Kollreider, Hartwig Fronthaler, and Josef Bigun. Non-intrusive liveness detection by face images. *Image and Vision Computing*, 27(3):233–244, 2009.
- [48] Wei Bao, Hong Li, Nan Li, and Wei Jiang. A liveness detection method for face recognition based on optical flow field. In *2009 International Conference on Image Analysis and Signal Processing*, pages 233–236. IEEE, 2009.
- [49] Klaus Kollreider, Hartwig Fronthaler, Maycel Isaac Faraj, and Josef Bigun. Real-time face detection and motion analysis with application in “liveness” assessment. *IEEE Transactions on Information Forensics and Security*, 2(3):548–558, 2007.
- [50] Xiaobai Li, Jukka Komulainen, Guoying Zhao, Pong-Chi Yuen, and Matti Pietikäinen. Generalized face anti-spoofing by detecting pulse from face videos. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 4244–4249. IEEE, 2016.
- [51] Ewa Magdalena Nowara, Ashutosh Sabharwal, and Ashok Veeraraghavan. Ppgsecure: Biometric presentation attack detection using photoplethysmograms. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*, pages 56–62. IEEE, 2017.
- [52] Siqi Liu, Pong C Yuen, Shengping Zhang, and Guoying Zhao. 3d mask face anti-spoofing with remote photoplethysmography. In *European Conference on Computer Vision*, pages 85–100. Springer, 2016.
- [53] Xiaoyang Tan, Yi Li, Jun Liu, and Lin Jiang. Face liveness detection from a single image with sparse low rank bilinear discriminative model. *Computer Vision—ECCV 2010*, pages 504–517, 2010.
- [54] Bruno Peixoto, Carolina Michelassi, and Anderson Rocha. Face liveness detection under bad illumination conditions. In *2011 18th IEEE International Conference on Image Processing*, pages 3557–3560. IEEE, 2011.
- [55] Jiamin Bai, Tian-Tsong Ng, Xinting Gao, and Yun-Qing Shi. Is physics-based liveness detection truly possible with a single image? In *Proceedings of 2010 IEEE International Symposium on Circuits and Systems*, pages 3425–3428. IEEE, 2010.
- [56] Jukka Määttä, Abdenour Hadid, and Matti Pietikäinen. Face spoofing detection from single images using micro-texture analysis. In *Biometrics (IJCB), 2011 international joint conference on*, pages 1–7. IEEE, 2011.
- [57] Neslihan Kose and Jean-Luc Dugelay. Countermeasure for the protection of face recognition systems against mask attacks. In *2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, pages 1–6. IEEE, 2013.

- [58] Neslihan Kose and Jean-Luc Dugelay. Shape and texture based countermeasure to protect face recognition systems against mask attacks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 111–116, 2013.
- [59] Nesli Erdogmus and Sébastien Marcel. Spoofing 2d face recognition systems with 3d masks. In *2013 International Conference of the BIOSIG Special Interest Group (BIOSIG)*, pages 1–8. IEEE, 2013.
- [60] Keyurkumar Patel, Hu Han, Anil K Jain, and Greg Ott. Live face video vs. spoof face video: Use of moiré patterns to detect replay video attacks. In *2015 International Conference on Biometrics (ICB)*, pages 98–105. IEEE, 2015.
- [61] Zinelabidine Boulkenafet, Jukka Komulainen, and Abdenour Hadid. Face spoofing detection using colour texture analysis. *IEEE Transactions on Information Forensics and Security*, 11(8):1818–1830, 2016.
- [62] Jukka Määttä, Abdenour Hadid, and Matti Pietikäinen. Face spoofing detection from single images using texture and local shape analysis. *IET biometrics*, 1(1):3–10, 2012.
- [63] Jianwei Yang, Zhen Lei, Shengcai Liao, and Stan Z Li. Face liveness detection with component dependent descriptor. In *2013 International Conference on Biometrics (ICB)*, pages 1–6. IEEE, 2013.
- [64] Jukka Komulainen, Abdenour Hadid, and Matti Pietikäinen. Context based face anti-spoofing. In *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*, pages 1–8, 2013.
- [65] Javier Galbally and Sébastien Marcel. Face anti-spoofing based on general image quality assessment. In *ICPR '14 Proceedings of the 2014 22nd International Conference on Pattern Recognition*, pages 1173–1178, 2014.
- [66] Javier Galbally, Sébastien Marcel, and Julian Fierrez. Image quality assessment for fake biometric detection: Application to iris, fingerprint, and face recognition. *IEEE transactions on image processing*, 23(2):710–724, 2013.
- [67] Keyurkumar Patel, Hu Han, and Anil K Jain. Secure face unlock: Spoof detection on smartphones. *IEEE transactions on information forensics and security*, 11(10):2268–2283, 2016.
- [68] Keyurkumar Patel, Hu Han, and Anil K Jain. Cross-database face antispoofing with robust feature representation. In *Chinese Conference on Biometric Recognition*, pages 611–619. Springer, 2016.
- [69] Lei Li, Xiaoyi Feng, Zinelabidine Boulkenafet, Zhaoqiang Xia, Mingming Li, and Abdenour Hadid. An original face anti-spoofing approach using partial convolutional neural network. In *2016 Sixth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pages 1–6. IEEE, 2016.
- [70] Amin Jourabloo, Yaojie Liu, and Xiaoming Liu. Face de-spoofing: Anti-spoofing via noise modeling. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 290–306, 2018.
- [71] Anjith George and Sébastien Marcel. Deep pixel-wise binary supervision for face presentation attack detection. In *2019 International Conference on Biometrics (ICB)*, pages 1–8. IEEE, 2019.
- [72] Tiago de Freitas Pereira and Anjos et al. Lbp-top based countermeasure against face spoofing attacks. In *ACCV*, pages 121–132. Springer, 2012.
- [73] Tiago de Freitas Pereira, Jukka Komulainen, André Anjos, José Mario De Martino, Abdenour Hadid, Matti Pietikäinen, and Sébastien Marcel. Face liveness detection using dynamic texture. *EURASIP Journal on Image and Video Processing*, 2014(1):2, 2014.
- [74] Tiago de Freitas Pereira, André Anjos, José Mario De Martino, and Sébastien Marcel. Can face anti-spoofing countermeasures work in a real world scenario? In *2013 international conference on biometrics (ICB)*, pages 1–8. IEEE, 2013.
- [75] Samarth Bharadwaj, Tejas I Dhamecha, and et al. Computationally efficient face spoofing detection with motion magnification. In *CVPR Workshops*, pages 105–110, 2013.
- [76] Allan da Silva Pinto, Helio Pedrini, William Schwartz, and Anderson Rocha. Video-based face spoofing detection through visual rhythm analysis. In *2012 25th SIBGRAPI Conference on Graphics, Patterns and Images*, pages 221–228. IEEE, 2012.
- [77] Allan Pinto, William Robson Schwartz, Helio Pedrini, and Anderson de Rezende Rocha. Using visual rhythms for detecting video-based facial spoof attacks. *IEEE Transactions on Information Forensics and Security*, 10(5):1025–1038, 2015.
- [78] Allan Pinto, Helio Pedrini, William Robson Schwartz, and Anderson Rocha. Face spoofing detection through visual codebooks of spectral temporal cubes. *IEEE Transactions on Image Processing*, 24(12):4726–4740, 2015.

- [79] Zhenqi Xu, Shan Li, and Weihong Deng. Learning temporal features using lstm-cnn architecture for face anti-spoofing. In *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, pages 141–145. IEEE, 2015.
- [80] Xiao Yang, Wenhan Luo, Linchao Bao, Yuan Gao, Dihong Gong, Shibao Zheng, Zhifeng Li, and Wei Liu. Face anti-spoofing: Model matters, so does data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3507–3516, 2019.
- [81] Tao Wang, Jianwei Yang, Zhen Lei, Shengcai Liao, and Stan Z Li. Face liveness detection using 3d structure recovered from a single camera. In *2013 international conference on biometrics (ICB)*, pages 1–6. IEEE, 2013.
- [82] Zezheng Wang, Chenxu Zhao, Yunxiao Qin, Qiusheng Zhou, Guojun Qi, Jun Wan, and Zhen Lei. Exploiting temporal and depth information for multi-frame face anti-spoofing. *arXiv preprint arXiv:1811.05118*, 2018.
- [83] Zitong Yu, Chenxu Zhao, Zezheng Wang, Yunxiao Qin, Zhuo Su, Xiaobai Li, Feng Zhou, and Guoying Zhao. Searching central difference convolutional networks for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5295–5305, 2020.
- [84] Gang Pan, Lin Sun, Zhaohui Wu, and Yueming Wang. Monocular camera-based face liveness detection by combining eyeblink and scene context. *Telecommunication Systems*, 47(3-4):215–225, 2011.
- [85] Litong Feng, Lai-Man Po, Yuming Li, Xuyuan Xu, Fang Yuan, Terence Chun-Ho Cheung, and Kwok-Wai Cheung. Integration of image quality and motion cues for face anti-spoofing: A neural network approach. *Journal of Visual Communication and Image Representation*, 38:451–460, 2016.
- [86] Olegs Nikisins, Amir Mohammadi, André Anjos, and Sébastien Marcel. On effectiveness of anomaly detection approaches against unseen presentation attacks in face anti-spoofing. In *2018 International Conference on Biometrics (ICB)*, pages 75–81. IEEE, 2018.
- [87] Shervin Rahimzadeh Arashloo, Josef Kittler, and William Christmas. An anomaly detection approach to face spoofing detection: A new formulation and evaluation protocol. *IEEE Access*, 5:13868–13882, 2017.
- [88] Yaojie Liu, Joel Stehouwer, Amin Jourabloo, and Xiaoming Liu. Deep tree learning for zero-shot face anti-spoofing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4680–4689, 2019.
- [89] R. Shao, X. Lan, J. Li, and P. Yuen. Multi-adversarial discriminative deep domain generalization for face presentation attack detection. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10015–10023, 2019.
- [90] JK Aggarwal and N Nandhakumar. On the computation of motion from sequences of images—a review. *Proceedings of the IEEE*, 76(8):917–935, 1988.
- [91] Ali Azarbayejani, Thad Starner, Bradley Horowitz, and Alex Pentland. Visually controlled graphics. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(6):602–605, 1993.
- [92] Josef Bigün, Goesta H. Granlund, and Johan Wiklund. Multidimensional orientation estimation with applications to texture analysis and optical flow. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (8):775–790, 1991.
- [93] Craig N Karson. Spontaneous eye-blink rates and dopaminergic systems. *Brain*, 106(3):643–653, 1983.
- [94] Lawrence R Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- [95] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition. CVPR 2001*, volume 1, pages I–I. IEEE, 2001.
- [96] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013.
- [97] E Howcroft. How faking videos became easy and why that’s so scary, 2018.
- [98] Robert Chesney and Danielle Citron. Deepfakes and the new disinformation war: The coming age of post-truth geopolitics. *Foreign Aff.*, 98:147, 2019.
- [99] Thanh Thi Nguyen, Cuong M Nguyen, Dung Tien Nguyen, Duc Thanh Nguyen, and Saeid Nahavandi. Deep learning for deepfakes creation and detection. *arXiv preprint arXiv:1909.11573*, 2019.

- [100] Stéphane Dupont and Juergen Luetttin. Audio-visual speech modeling for continuous speech recognition. *IEEE transactions on multimedia*, 2(3):141–151, 2000.
- [101] Deepfakes web. <https://deepfakesweb.com/>. (accessed on 15/09/2020).
- [102] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017.
- [103] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1947–1962, 2018.
- [104] Faceapp. <https://www.faceapp.com/>. (accessed on 15/09/2020).
- [105] Fakeapp. <https://www.fakeapp.org/>. (accessed on 15/09/2020).
- [106] Steven Fernandes, Sunny Raj, Eddy Ortiz, Iustina Vintila, Margaret Salter, Gordana Urosevic, and Sumit Jha. Predicting heart rate variations of deepfake videos using neural ode. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [107] Umur Aybars Ciftci, Ilke Demir, and Lijun Yin. Fakecatcher: Detection of synthetic portrait videos using biological signals. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [108] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5001–5010, 2020.
- [109] Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.
- [110] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in neural information processing systems*, pages 6571–6583, 2018.
- [111] Michael Oren and Shree K Nayar. Generalization of the lambertian model and implications for machine vision. *International Journal of Computer Vision*, 14(3):227–251, 1995.
- [112] Xiaoyang Tan and Bill Triggs. Enhanced local texture feature sets for face recognition under difficult lighting conditions. *IEEE transactions on image processing*, 19(6):1635–1650, 2010.
- [113] Karel Zuiderveld. Contrast limited adaptive histogram equalization. *Graphics gems*, pages 474–485, 1994.
- [114] Robby T Tan and Katsushi Ikeuchi. Separating reflection components of textured surfaces using a single image. In *Digitally Archiving Cultural Objects*, pages 353–384. Springer, 2008.
- [115] Timo Ojala, Matti Pietikainen, and Topi Maenpaa. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on pattern analysis and machine intelligence*, 24(7):971–987, 2002.
- [116] Timo Ahonen, Abdenour Hadid, and Matti Pietikainen. Face description with local binary patterns: Application to face recognition. *IEEE transactions on pattern analysis and machine intelligence*, 28(12):2037–2041, 2006.
- [117] Lun Zhang, Rufeng Chu, Shiming Xiang, Shengcai Liao, and Stan Z Li. Face detection based on multi-block lbp representation. In *International conference on biometrics*, pages 11–18. Springer, 2007.
- [118] Guoying Zhao and Matti Pietikainen. Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE transactions on pattern analysis and machine intelligence*, 29(6):915–928, 2007.
- [119] Bangalore S Manjunath and Wei-Ying Ma. Texture features for browsing and retrieval of image data. *IEEE Transactions on pattern analysis and machine intelligence*, 18(8):837–842, 1996.
- [120] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, volume 1, pages 886–893. IEEE, 2005.
- [121] Andrea Vedaldi and Andrew Zisserman. Efficient additive kernels via explicit feature maps. *IEEE transactions on pattern analysis and machine intelligence*, 34(3):480–492, 2012.
- [122] Vittorio Ferrari, Manuel Marin-Jimenez, and Andrew Zisserman. Progressive search space reduction for human pose estimation. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008.

- [123] Ville Ojansivu and Janne Heikkilä. Blur insensitive texture classification using local phase quantization. In *International conference on image and signal processing*, pages 236–243. Springer, 2008.
- [124] Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 2169–2178. IEEE, 2006.
- [125] Xinting Gao, Tian-Tsong Ng, Bo Qiu, and Shih-Fu Chang. Single-view recaptured image detection based on physics-based features. In *2010 IEEE International Conference on Multimedia and Expo*, pages 1469–1474. IEEE, 2010.
- [126] Frederique Crete, Thierry Dolmiere, Patricia Ladret, and Marina Nicolas. The blur effect: perception and estimation with a new no-reference perceptual blur metric. In *Human vision and electronic imaging XII*, volume 6492, page 64920I. International Society for Optics and Photonics, 2007.
- [127] Pina Marziliano, Frederic Dufaux, Stefan Winkler, and Touradj Ebrahimi. A no-reference perceptual blur metric. In *Proceedings. International Conference on Image Processing*, volume 3, pages III–III. IEEE, 2002.
- [128] Yuanhao Chen, Zhiwei Li, Mingjing Li, and Wei-Ying Ma. Automatic classification of photographs and graphics. In *2006 IEEE International Conference on Multimedia and Expo*, pages 973–976. IEEE, 2006.
- [129] Anil Jain, Karthik Nandakumar, and Arun Ross. Score normalization in multimodal biometric systems. *Pattern recognition*, 38(12):2270–2285, 2005.
- [130] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [131] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014.
- [132] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [133] Hao-Yu Wu, Michael Rubinstein, Eugene Shih, John Guttag, Frédo Durand, and William Freeman. Eulerian video magnification for revealing subtle changes in the world. *ACM transactions on graphics (TOG)*, 31(4):1–8, 2012.
- [134] Rizwan Chaudhry, Avinash Ravichandran, Gregory Hager, and René Vidal. Histograms of oriented optical flow and binet-cauchy kernels on nonlinear dynamical systems for the recognition of human actions. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1932–1939. IEEE, 2009.
- [135] Seong Soo Chun, Hyeokman Kim, Kim Jung-Rim, Sangwook Oh, and Sanghoon Sull. Fast text caption localization on video using visual rhythm. In *International Conference on Advances in Visual Information Systems*, pages 259–268. Springer, 2002.
- [136] Silvio Jamil Ferzoli Guimar, Michel Couprie, Neucimar Jerónimo Leite, et al. A method for cut detection based on visual rhythm. In *Proceedings XIV Brazilian Symposium on Computer Graphics and Image Processing*, pages 297–304. IEEE, 2001.
- [137] Robert M. Haralick, K. Sam Shanmugam, and Its'hak Dinstein. Textural features for image classification. *IEEE Trans. Syst. Man Cybern.*, 3:610–621, 1973.
- [138] Josef Sivic and Andrew Zisserman. Video google: A text retrieval approach to object matching in videos. In *null*, page 1470. IEEE, 2003.
- [139] Peter J Schmid, Larry Li, Matthew P Juniper, and O Pust. Applications of the dynamic mode decomposition. *Theoretical and Computational Fluid Dynamics*, 25(1-4):249–259, 2011.
- [140] Jason M Saragih, Simon Lucey, and Jeffrey F Cohn. Deformable model fitting by regularized landmark mean-shift. *International journal of computer vision*, 91(2):200–215, 2011.
- [141] Richard I Hartley and Peter Sturm. Triangulation. *Computer vision and image understanding*, 68(2):146–157, 1997.
- [142] Amin Jourabloo and Xiaoming Liu. Large-pose face alignment via cnn-based dense 3d model fitting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4188–4196, 2016.
- [143] Amin Jourabloo and Xiaoming Liu. Pose-invariant face alignment via cnn-based dense 3d model fitting. *International Journal of Computer Vision*, 124(2):187–203, 2017.

- [144] Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 534–551, 2018.
- [145] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [146] Volker Blanz and Thomas Vetter. Face recognition based on fitting a 3d morphable model. *IEEE Transactions on pattern analysis and machine intelligence*, 25(9):1063–1074, 2003.
- [147] Takashi Matsumoto. Graphics system shadow generation using a depth buffer, 1991. US Patent 5,043,922.
- [148] André Anjos and Sébastien Marcel. Counter-measures to photo attacks in face recognition: a public database and a baseline. In *2011 international joint conference on Biometrics (IJCB)*, pages 1–7. IEEE, 2011.
- [149] Diego A Socolinsky, Andrea Selinger, and Joshua D Neuheisel. Face recognition with visible and thermal infrared imagery. *Computer vision and image understanding*, 91(1-2):72–114, 2003.
- [150] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.
- [151] Murali Mohan Chakka, Andre Anjos, Sebastien Marcel, Roberto Tronci, Daniele Muntoni, Gianluca Fadda, Maurizio Pili, Nicola Sirena, Gabriele Murgia, Marco Ristori, et al. Competition on counter measures to 2-d facial spoofing attacks. In *2011 International Joint Conference on Biometrics (IJCB)*, pages 1–6. IEEE, 2011.
- [152] Ivana Chingovska, Jinwei Yang, Zhen Lei, Dong Yi, Stan Z Li, Olga Kahm, Christian Glaser, Naser Damer, Arjan Kuijper, Alexander Nouak, et al. The 2nd competition on counter measures to 2d face spoofing attacks. In *2013 International Conference on Biometrics (ICB)*, pages 1–6. IEEE, 2013.
- [153] Junjie Yan, Zhiwei Zhang, Zhen Lei, Dong Yi, and Stan Z Li. Face liveness detection by exploring multiple scenic clues. In *2012 12th International Conference on Control Automation Robotics & Vision (ICARCV)*, pages 188–193. IEEE, 2012.
- [154] Yasemin Altun, Ioannis Tsochantaridis, and Thomas Hofmann. Hidden markov support vector machines. In *Proceedings of the 20th international conference on machine learning (ICML-03)*, pages 3–10, 2003.
- [155] Glenn Easley, Demetrio Labate, and Wang-Q Lim. Sparse directional image representations using the discrete shearlet transform. *Applied and Computational Harmonic Analysis*, 25(1):25–46, 2008.
- [156] Yuming Li, Lai-Man Po, Xuyuan Xu, and Litong Feng. No-reference image quality assessment using statistical characterization in the shearlet domain. *Signal Processing: Image Communication*, 29(7):748–759, 2014.
- [157] Ce Liu et al. *Beyond pixels: exploring new representations and applications for motion analysis*. PhD thesis, Massachusetts Institute of Technology, 2009.
- [158] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine*, 29(6):82–97, 2012.
- [159] Yann LeCun, Bernhard Boser, and et al. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [160] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- [161] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [162] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261*, 2016.
- [163] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Smash: one-shot model architecture search through hypernetworks. *arXiv preprint arXiv:1708.05344*, 2017.
- [164] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.

- [165] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018.
- [166] Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Jie Tan, Quoc Le, and Alex Kurakin. Large-scale evolution of image classifiers. *arXiv preprint arXiv:1703.01041*, 2017.
- [167] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 4780–4789, 2019.
- [168] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018.
- [169] Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong. Pc-darts: Partial channel connections for memory-efficient differentiable architecture search. *arXiv preprint arXiv:1907.05737*, 2019.
- [170] Ning Zhu. Neural architecture search for deep face recognition. *arXiv preprint arXiv:1904.09523*, 2019.
- [171] Wei Peng, Xiaopeng Hong, and Guoying Zhao. Video action recognition via neural architecture searching. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 11–15. IEEE, 2019.
- [172] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3750–3759, 2019.
- [173] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7036–7045, 2019.
- [174] Yiheng Zhang, Zhaofan Qiu, Jingen Liu, Ting Yao, Dong Liu, and Tao Mei. Customizable architecture search for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11641–11650, 2019.
- [175] Zinelabdine Boulkenafet, Jukka Komulainen, Zahid Akhtar, Azeddine Benlamoudi, Djamel Samai, Salah Eddine Bekhouche, Abdelkrim Ouafi, Fadi Dornaika, Abdelmalik Taleb-Ahmed, Le Qin, et al. A competition on generalized software-based face presentation attack detection in mobile scenarios. In *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pages 688–696. IEEE, 2017.
- [176] David Martinus Johannes Tax. One-class classification: Concept learning in the absence of counter-examples. 2002.
- [177] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010.
- [178] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528. IEEE, 2011.
- [179] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550, 2017.
- [180] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5400–5409, 2018.
- [181] Zinelabinde Boulkenafet, Jukka Komulainen, and et al. Oulu-npu: A mobile face presentation attack database with real-world variations. In *International Conference on Automatic Face & Gesture Recognition*, pages 612–618. IEEE, 2017.
- [182] Shifeng Zhang, Xiaobo Wang, Ajian Liu, Chenxu Zhao, Jun Wan, Sergio Escalera, Hailin Shi, Zezheng Wang, and Stan Z Li. A dataset and benchmark for large-scale multi-modal face anti-spoofing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 919–928, 2019.
- [183] Dayong Wang, Steven CH Hoi, Ying He, Jianke Zhu, Tao Mei, and Jiebo Luo. Retrieval-based face annotation by weak label regularized local coordinate coding. *IEEE transactions on pattern analysis and machine intelligence*, 36(3):550–563, 2013.
- [184] Ira Kemelmacher-Shlizerman, Steven M Seitz, Daniel Miller, and Evan Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4873–4882, 2016.

- [185] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [186] Paul Viola and Michael J Jones. Robust real-time face detection. *International journal of computer vision*, 57(2):137–154, 2004.
- [187] Xiaoyang Tan, Fengyi Song, Zhi-Hua Zhou, and Songcan Chen. Enhanced pictorial structures for precise eye localization under uncontrolled conditions. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1621–1628. IEEE, 2009.
- [188] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pages 2818–2826, 2016.
- [189] André Anjos, Murali Mohan Chakka, and Sébastien Marcel. Motion-based counter-measures to photo attacks in face recognition. *IET biometrics*, 3(3):147–158, 2013.
- [190] Cunjian Chen, Antitza Dantcheva, and Arun Ross. Automatic facial makeup detection with application in face recognition. In *2013 international conference on biometrics (ICB)*, pages 1–8. IEEE, 2013.
- [191] Huiwen Chang, Jingwan Lu, Fisher Yu, and Adam Finkelstein. Pairedcyclegan: Asymmetric style transfer for applying and removing makeup. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 40–48, 2018.