

Face Anti-Spoofing with Human Material Perception

Zitong Yu¹, Xiaobai Li¹, Xuesong Niu^{2,3}, Jingang Shi^{4,1}, Guoying Zhao^{1*}

¹Center for Machine Vision and Signal Analysis, University of Oulu, Finland

²Key Laboratory of Intelligent Information Processing of Chinese Academy of Sciences (CAS), Institute of Computing Technology, CAS, China

³University of Chinese Academy of Sciences, China

⁴School of Software Engineering, Xian Jiaotong University, China

{zitong.yu, xiaobai.li, guoying.zhao}@oulu.fi, {xuesong.niu}@vip1.ict.ac.cn

Abstract. Face anti-spoofing (FAS) plays a vital role in securing the face recognition systems from presentation attacks. Most existing FAS methods capture various cues (e.g., texture, depth and reflection) to distinguish the live faces from the spoofing faces. All these cues are based on the discrepancy among physical materials (e.g., skin, glass, paper and silicone). In this paper we rephrase face anti-spoofing as a material recognition problem and combine it with classical human material perception [1], intending to extract discriminative and robust features for FAS. To this end, we propose the Bilateral Convolutional Networks (BCN), which is able to capture intrinsic material-based patterns via aggregating multi-level bilateral macro- and micro- information. Furthermore, Multi-level Feature Refinement Module (MFRM) and multi-head supervision are utilized to learn more robust features. Comprehensive experiments are performed on six benchmark datasets, and the proposed method achieves superior performance on both intra- and cross-dataset testings. One highlight is that we achieve overall $11.3 \pm 9.5\%$ EER for cross-type testing in SiW-M dataset, which significantly outperforms previous results. We hope this work will facilitate future cooperation between FAS and material communities.

1 Introduction

In recent years, face recognition has been widely used in various interactive and payment scene due to its high accuracy and convenience. However, such biometric system is vulnerable to presentation attacks (PAs). Typical examples of physical presentation attacks include print, video replay, 3D masks and makeup. In order to detect such PAs and secure the face recognition system, face anti-spoofing (FAS) has attracted more attention from both academia and industry.

In the past decade, several hand-crafted feature based [2,3,4,5,6,7] and deep learning based [8,9,10,11] methods have been proposed for presentation attack detection (PAD). On one hand, the classical hand-crafted descriptors leverage local relationship among the neighbours as the discriminative features, which

* Corresponding author

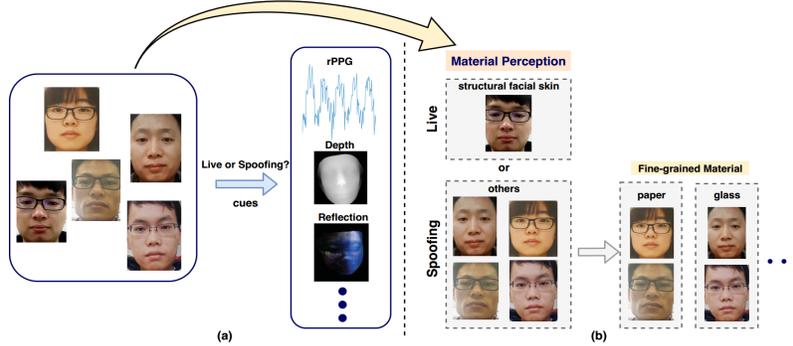


Fig. 1. (a) Face anti-spoofing can be regarded as a binary classification (live or spoofing) problem, which relies on the intrinsic cues such as rPPG, depth, reflection and so on. (b) Face anti-spoofing can be also treated as a material perception problem.

is robust for describing the detailed invariant information (e.g., color texture, moiré pattern and noise artifacts) between the live and spoofing faces. On the other hand, due to the stacked convolution operations with nonlinear activation, the convolutional neural networks (CNN) hold strong representation abilities to distinguish the bona fide from PA. However, most existing CNN and hand-crafted features are designed for universal image recognition tasks, which might not represent fine-grained spoofing patterns in FAS task.

According to the known intrinsic cues in face anti-spoofing task, many state-of-the-art methods introduced task-oriented priori knowledge for feature representation. As shown in Fig. 1(a), there are three famous human-defined cues (i.e., rPPG, depth and reflection) for FAS task. Firstly, frequency distribution dissimilarity of rPPG signals [12,13,14,15] recovered from live skin surface and spoofing face can be utilized as there are no or weaker blood volume changes in spoofing faces. Secondly, structural facial depth difference between live and spoofing faces [14,16] can be adopted as significant cue as most spoofing faces are broadcasted in plane presentation attack instruments (PAIs). Thirdly, reflectance difference [17,18] is also one kind of reliable cues as human facial skin and spoofing surfaces react differently to changes in illumination. Despite the human-defined cues are helpful to enhance the modeling capability respectively, it is still difficult to describe intrinsic and robust features for FAS task.

An interesting and essential question for FAS task is how human beings differentiate live or spoofing faces, and what can be learned by machine intelligent systems? In real-world cases, spoofing faces are always broadcasted by physical spoofing carriers (e.g., paper, glass screen and resin mask), which have obvious material properties difference with human facial skin. Such difference can be explicitly described as human-defined cues (e.g., rPPG, depth and reflection) or implicitly learned according to the material property uniqueness of structural live facial skin. Therefore, as illustrated in Fig. 1(b), we assume that discrepancy of the structural materials between human facial skin

and physical spoofing carriers are the essence of distinguishing live faces from spoofing ones.

Motivated by the discussions above, we rephrase face anti-spoofing task as structural material recognition problem and our goal is to learn intrinsic and robust features for distinguishing structural facial skin material from the others (i.e., materials for physical spoofing carriers). According to the study inspired by classical human material perception [1], bilateral filtering plays a vital role in representing macro- and micro- cues for various materials. In this paper, we integrate traditional bilateral filtering operator into the state-of-the-art FAS deep learning framework, intending to help networks to learn more intrinsic material-based patterns. Our contributions include:

- We design novel Bilateral Convolutional Networks (BCN), which is able to capture intrinsic material-based patterns via aggregating multi-level bilateral macro- and micro- information.
- We propose to use Multi-level Feature Refinement Module (MFRM) and material based multi-head supervision to further boost the performance of BCN. The former one refines the multi-scale features via reassembling weights of local neighborhood while the latter forces the network to learn robust shared features for multi-head auxiliary tasks.
- Our proposed method achieves outstanding performance on six benchmark datasets with both intra- and cross-dataset testing protocols. We also conduct fine-grained material recognition experiments on SiW-M dataset to validate the effectiveness of our proposed method.

2 Related Work

2.1 Face Anti-Spoofing

Traditional face anti-spoofing methods usually extract hand-crafted features from the facial images to capture the spoofing patterns. Several classical local descriptors such as LBP [2,4], SIFT [7], SURF [19], HOG [5] and DoG [6] are utilized to extract frame level features while video level methods usually capture dynamic cues like dynamic texture [20], micro-motion [21] and eye blinking [22]. More recently, a few deep learning based methods are proposed for FAS task. Some frame-level CNN methods [23,24,25,26,11,27] are supervised by binary scalars or pixel-wise binary maps. In contrast, auxiliary depth [10,16,14] and reflection [28] supervisions are introduced to learn detailed cues effectively. In order to learn generalized features for unseen attacks and environment, few-shot learning [8], zero-shot learning [8,29] and domain generalization [30,31,32] are introduced for FAS task. Meanwhile, several video-level CNN methods are presented to exploit the dynamic spatio-temporal [9,33,34,35] or rPPG [13,14,12,36,37,38] features for PAD. Despite introducing task-oriented priori knowledge (e.g., auxiliary depth, reflection and rPPG), deep learning based methods are still difficult to extract rich intrinsic features among live faces and various kinds of PAs.

2.2 Human and Machine Material Perception

Our world consists of not only objects and scenes but also of materials of various kinds. The perception of materials by humans usually focuses on optical and

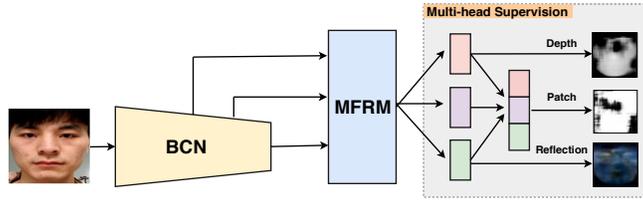


Fig. 2. The overall framework consists of Bilateral Convolutional Networks (BCN), Multi-level Feature Refinement Module (MFRM) and multi-head supervision.

mechanical properties. Maloney and Brainard [39] demonstrates the research concerns about perception of material surface properties other than color and lightness, such as gloss or roughness. Fleming [40] proposes statistical appearance models to describe visual perception of materials. Nishida [41] presents that material perception is visual estimation of optical modulation of image statistics. Inspired by human material perception, several machine intelligent methods are designed for material classification. Techniques derived from the domain of texture analysis can be adopted for material recognition by machines [42]. Varma and Zisserman [43] utilizes joint distribution of intensity values over image patch exemplars for material classification under unknown viewpoint and illumination. Sharan et al. [1] uses bilateral based low and mid-level image features for material recognition. Aiming to keep the details of features, deep dilated convolutional network is used for material perception [44].

In terms of vision applications, concepts of human material perception have been developed into image quality assessment [45] and video quality assessment [46]. For face anti-spoofing task, few works [17,18,47] consider discrepant surface reflectance properties of live or spoofing faces. However, only considering surface reflectance properties is not always reliable for material perception [1]. In order to learn more generalized material-based features for FAS, we combine the state-of-the-art FAS methods with classical human material perception[1].

3 Methodology

In this section, we first introduce the Bilateral Convolutional Networks (BCN) in Section 3.1, then present Multi-level Feature Refinement Module (MFRM) in Section 3.2, and at last introduce the material based multi-head supervision for face anti-spoofing in Section 3.3. The overall framework is shown in Fig. 2.

3.1 Bilateral Convolutional Networks

Inspired by classical material perception [1] that utilizes bilateral filtering [48] for exacting subsequent macro- and micro- features, we try to adopt bilateral filtering technique for FAS task. The main issue is that in [1], several hand-crafted features are designed after bilateral filtering, which limits the feature representation capacity. In this subsection, we propose two solutions to integrate bilateral filtering with the state-of-the-art deep networks for FAS task.

Bilateral Filtering. The first solution is straightforward: The bilateral filtered frames are taken as network inputs instead of the original RGB frames.

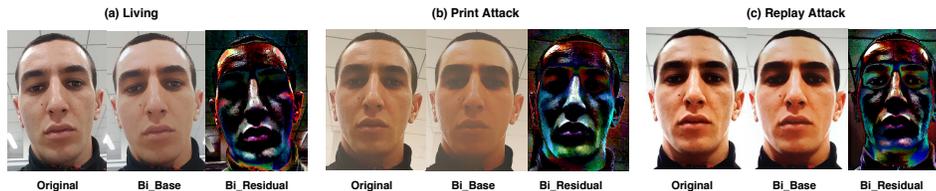


Fig. 3. Samples visualization for (a) live faces, (b) print attack, and (c) replay attack. ‘Bi_Base’ denotes frames after bilateral filtering while ‘Bi_Residual’ denotes residual result between original and bilateral filtered frames respectively. The intensity values of bilateral residual images are enlarged by four times for better visual effects.

The bilateral filter is utilized to smooth the original frame while preserving its main edges. Each pixel is a weighted mean of its neighbors where the weights decrease with the distance in space and with the intensity difference. With Gaussian function $g_{\sigma}(x) = \exp(-x^2/\sigma^2)$, the bilateral filter of image I at pixel p is defined by:

$$Bi_Base(I)_p = \frac{1}{k} \sum_{q \in I} g_{\sigma_s}(\|p - q\|) g_{\sigma_r}(|I_p - I_q|) I_q, \quad (1)$$

$$with : \quad k = \sum_{q \in I} g_{\sigma_s}(\|p - q\|) g_{\sigma_r}(|I_p - I_q|),$$

where σ_s and σ_r control the influence of spatial neighborhood distance and intensity difference respectively, and k normalizes the weights. Give the input image I , bilateral filter is able to create a two-scale decomposition [49] where the output of the filter produces a large-scale base image $Bi_Base(I)$ and the residual detail image $Bi_Residual(I)$ can be obtained by $Bi_Residual(I) = I - Bi_Base(I)$. We use the fast approximation version¹ of the bilateral filter [50] with default parameters for implementation.

Typical samples before and after bilateral filtering are visualized in Figure 3. There are obvious differences in bilateral base and residual images between live and spoofing faces despite their similarities in the original RGB images. As shown in Fig. 3(b) ‘Bi_Base’, the print attack face made of paper material is rougher and less glossy. Moreover, it can be seen from Fig. 3(b)(c) ‘Bi_Residual’ that the high-frequency activation in eyes and eyebrow region is stronger, which might be caused by discrepant surface reflectance properties among materials (e.g., skin, paper and glass). These visual evidences are consistent with classical human material perception [1] that macro- cues from bilateral base and micro- cues from bilateral residual are helpful for material perception.

In this paper, Auxiliary(Depth) [14] is chosen as our baseline deep model. The bilateral filtered (i.e., bilateral base and residual) images can forward the baseline model directly and predict the corresponding results. The ablation study of different kinds of inputs will be discussed in Section 4.3.

Deep Bilateral Networks. The drawbacks of the above-mentioned solution are mainly of two folds: 1) directly replacing original inputs with bilateral

¹ <http://people.csail.mit.edu/jiawen/software/bilateralFilter-1.0.m>

images might lead to information loss, which limits the feature representation capability for neural networks, and 2) it is an inefficient way to learn multi-level bilateral features as the bilateral filter is only adopted in the input space. Aiming to overcome these drawbacks, we propose a novel method called Bilateral Convolutional Networks (BCN) to integrate traditional bilateral filtering with deep networks properly.

In order to filter the deep features instead of original images, the deep bilateral operator (DBO) is introduced. Mimicking the process of gray-scale or color image filtering, given the deep feature maps $\mathcal{F} \in \mathbb{R}^{H \times W \times C}$ with height H , width W and C channels, channel-wise deep bilateral filtering is operated. Considering the small spatial distance for the widely used convolution with 3×3 kernel, the distance decay term in Eqn. (1) can be removed (see *Appendix A* for corresponding ablation study), which is more efficient and lightweight when operating in deep hidden space. Hence deep bilateral operator for each channel of \mathcal{F} can be formulated as

$$DBO(\mathcal{F})_p = \frac{1}{k} \sum_{q \in \mathcal{F}} g_{\sigma_r}(|\mathcal{F}_p - \mathcal{F}_q|) \mathcal{F}_q, \quad (2)$$

$$with : \quad k = \sum_{q \in \mathcal{F}} g_{\sigma_r}(|\mathcal{F}_p - \mathcal{F}_q|).$$

Now performing DBO for features in different levels, it is easy to obtain multi-level bilateral base features. Nevertheless, how to get multi-level bilateral residual features is still unknown. As our goal is to represent aggregated bilateral base and residual features \mathcal{F}_{Bi} , inspired by residual learning in ResNet [51], bilateral residual features $\mathcal{F}_{Residual}$ can be learned dynamically via shortcut connecting with bilateral base features \mathcal{F}_{Base} , i.e., $\mathcal{F}_{Residual} = \mathcal{F}_{Bi} - \mathcal{F}_{Base}$. The architecture of the proposed Bilateral Convolutional Networks is illustrated in Figure 4. As ‘BilateralConvBlock’ and ‘ConvBlock’ have same convolutional structure but unshared parameters, it is possible to learn \mathcal{F}_{Base} and $\mathcal{F}_{Residual}$ from ‘BilateralConvBlock’ and ‘ConvBlock’ respectively. Compared with baseline model Auxiliary(Depth) [14] without ‘BilateralConvBlock’, BCN is able to learn more intrinsic features via aggregating multi-level bilateral macro- and micro- information.

3.2 Multi-level Feature Refinement Module

In the baseline model Auxiliary(Depth) [14], multi-level features are concatenated directly for subsequent head supervision. We argue that such coarse features are not optimal for fine-grained material-based FAS task. Hence Multi-level Feature Refinement Module (MFRM) is introduced after BCN (see Fig. 2), which aims to refine and fuse the coarse low-mid-high level features from BCN via context-aware feature reassembling.

As illustrated in Fig. 5, features \mathcal{F} from low-mid-high levels are refined via reassembling features with local context-aware weights. Unlike [52] which aims for feature upsampling, here we focus on the general multi-level feature refinement. The refined features \mathcal{F}' can be formulated as

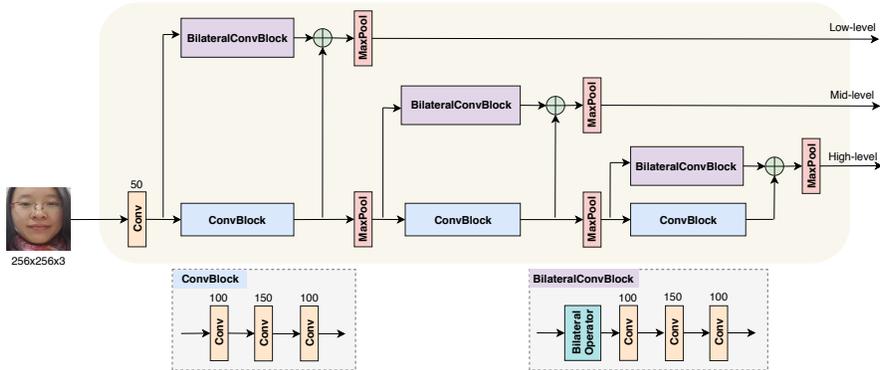


Fig. 4. The proposed BCN architecture. The number of filters are shown on top of each convolutional layer, the size of all filters is 3×3 with stride 1 for convolutional and 2 for pooling layers. Each output from ‘ConvBlock’ and ‘BilateralConvBlock’ in the same level will be operated with element-wise addition.

$$\mathcal{F}'_{level} = \mathcal{F}_{level} \otimes \mathcal{N}(\psi(\phi(\mathcal{F}_{level}))), \quad level \in \{low, mid, high\}, \quad (3)$$

where ϕ , ψ , \mathcal{N} and \otimes represent channel compressor, content encoder, kernel normalizer and refinement operator, respectively. The channel compressor adopts a 1×1 convolution layer to compress the input feature channel from C to C' , making the refinement module more efficient. The content encoder utilizes a convolution layer of kernel size 5×5 to generate refinement kernels based on the content of input features \mathcal{F} , and then each $K \times K$ refinement kernel is normalized with a softmax function spatially in kernel normalizer. Given the location $l = (i, j)$, channel c and corresponding normalized refinement kernel \mathcal{W}_l , the output refined features \mathcal{F}' are expressed as

$$\mathcal{F}'_{(i,j,c)} = \sum_{n=-r}^r \sum_{m=-r}^r \mathcal{W}_{l(n,m)} \cdot \mathcal{F}_{(i+n,j+m,c)}, \quad with \quad r = \lfloor K/2 \rfloor. \quad (4)$$

In essence, MFRM exploits the semantic and contextual information to reallocate the contributions of the local neighbors, which is possible to obtain more intrinsic features. For instance, our module is able to refine the discriminative cues (e.g. moiré pattern) from salient regions according to their local context. We also compare the refinement method with other classical feature attention based methods such as spatial attention [53], channel attention [54] and non-local attention [55] in Sec. 4.3.

3.3 Material based Multi-head Supervision

As material categorization is complex and depends on various fine-grained cues (e.g., surface shape, reflectance and texture), it is impossible to learn robust and intrinsic material-based patterns via only one simple supervision (e.g., softmax binary loss and depth regression loss). For the sake of learning intrinsic material-based features, material based multi-head supervision is proposed. As shown in Fig. 2, three-head supervision is utilized to guide the multi-level fused features

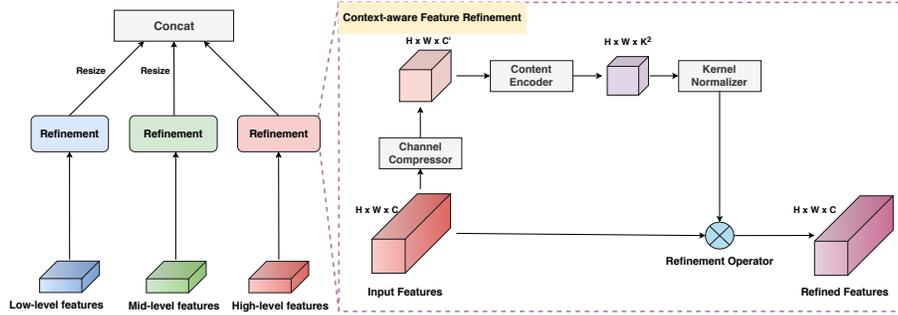


Fig. 5. Multi-level Feature Refinement Module.

from MFRM: 1) depth-head supervision, intending to force the model to learn structural surface shape information; 2) reflection-head supervision, helping networks to learn surface reflectance property; and 3) patch-head supervision, guiding to learn fine-grained surface texture cues. The detailed network structure can be found in *Appendix B*.

Loss Function. Appropriate loss functions should be designed to supervise the network training. Given an input face image I , the network predicts the depth map D_{pre} , reflection map R_{pre} and patch map P_{pre} . Then the loss functions can be formulated as:

$$\mathcal{L}_{depth} = \frac{1}{H \times W} \sum_{i \in H, j \in W} \|D_{pre(i,j)} - D_{gt(i,j)}\|_2^2, \quad (5)$$

$$\mathcal{L}_{reflection} = \frac{1}{H \times W \times C} \sum_{i \in H, j \in W, c \in C} \|R_{pre(i,j,c)} - R_{gt(i,j,c)}\|_2^2, \quad (6)$$

$$\mathcal{L}_{patch} = \frac{1}{H \times W} \sum_{i \in H, j \in W} -(P_{gt(i,j)} \log(P_{pre(i,j)}) + (1 - P_{gt(i,j)}) \log(1 - P_{pre(i,j)})), \quad (7)$$

where D_{gt} , R_{gt} and P_{gt} denote ground truth depth map, reflection map and patch map respectively. Finally, the overall loss function is $\mathcal{L}_{overall} = \mathcal{L}_{depth} + \mathcal{L}_{reflection} + \mathcal{L}_{patch}$.

Ground Truth Generation. Dense face alignment PRNet [56] is adopted to estimate facial 3D shapes and generate the facial depth maps with size 32×32 . The reflection maps are estimated by the state-of-the-art reflection estimation network [57] and then face regions are cropped from reflection maps to avoid being overfitted to backgrounds. The generated reflection maps have size with $32 \times 32 \times 3$ (3 channels for RGB). More details and samples can be found in [33,28]. To distinguish live faces from spoofing faces, at the training stage, we normalize live depth maps and spoofing reflection maps in a range of $[0, 1]$, while setting spoofing depth maps and live reflection maps to 0, which is similar to [14,28]. The patch maps are generated simply by downsampling original images and filling each patch position with corresponding binary label (i.e., live 1 and spoofing 0). The generated binary patch maps keep size with 32×32 .

4 Experiments

In this section, comprehensive experiments are performed to evaluate our method. We will sequentially describe the employed datasets & metrics (Sec. 4.1), implementation details (Sec. 4.2), results (Sec. 4.3 - 4.5) and analysis (Sec. 4.6).

4.1 Datasets and Metrics

Six databases including OULU-NPU [58], SiW [14], CASIA-MFSD [59], Replay-Attack [60], MSU-MFSD [61] and SiW-M [29] are used in our experiments. OULU-NPU and SiW are large-scale high-resolution databases, containing four and three protocols to validate the generalization (e.g., unseen illumination and attack medium) of models respectively, which are utilized for intra testing. CASIA-MFSD, Replay-Attack and MSU-MFSD are databases which contain low-resolution videos, and are used for cross testing. SiW-M is designed for fine-grained material recognition and cross-type testing for unseen attacks as there are rich attacks types (totally 13 types) inside.

Performance Metrics. In OULU-NPU and SiW dataset, we follow the original protocols and metrics, i.e., Attack Presentation Classification Error Rate (APCER), Bona Fide Presentation Classification Error Rate (BPCER), and ACER [62] for a fair comparison. Half Total Error Rate (HTER) is adopted in the cross testing between CASIA-MFSD and Replay-Attack. Area Under Curve (AUC) is utilized for intra-database cross-type test on CASIA-MFSD, Replay-Attack and MSU-MFSD. For the cross-type test on SiW-M, APCER, BPCER, ACER and Equal Error Rate (EER) are employed.

4.2 Implementation Details

Our proposed method is implemented with Pytorch. The default settings $\sigma_r = 1.0$ and $C' = 20, K = 5$ are adopted for BCN and MFRM, respectively. In the training stage, models are trained with Adam optimizer and the initial learning rate (lr) and weight decay (wd) are $1e-4$ and $5e-5$, respectively. We train models with maximum 1300 epochs while lr halves every 500 epochs. The batch size is 7 on a Nvidia P100 GPU. In the testing stage, we calculate the mean value of the predicted depth map \mathcal{D}_{test} , reflection map \mathcal{R}_{test} and patch map \mathcal{P}_{test} as the final score \mathcal{S}_{test} :

$$\mathcal{S}_{test} = mean(\mathcal{D}_{test}) + mean(1 - \mathcal{R}_{test}) + mean(\mathcal{P}_{test}). \quad (8)$$

4.3 Ablation Study

In this subsection, all ablation studies are conducted on Protocol-1 (different illumination condition and location between train and test sets) of OULU-NPU [58] to explore the details of our proposed BCN, MFRM and multi-head supervision.

Impact of σ_r and Bilateral Operator in BCN. According to Eqn. (2) and Gaussian function $g_{\sigma_r}(x) = exp(-x^2/\sigma_r^2)$, σ_r controls the strength of neighbor feature differences, i.e., the higher σ_r , the more contributions are given to the neighbor with large differences. As illustrated in Fig. 6(a), the best performance

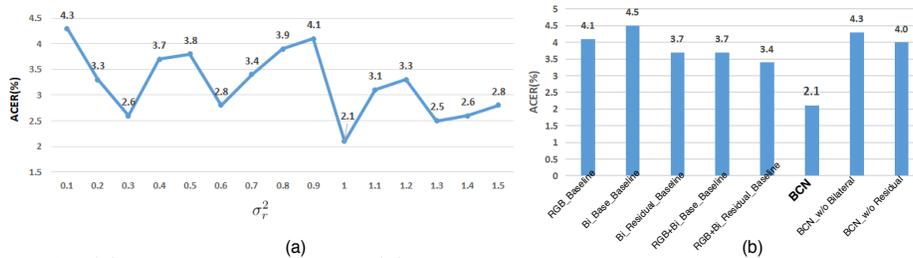


Fig. 6. (a) Impact of σ_r in BCN. (b) Comparison among various kinds of inputs for baseline model Auxiliary(Depth) [14] and BCN. Lower ACER, better performance.

Table 1. Results of network composition and supervision.

Model	ACER(%)
D (Baseline)	4.1
D+R	2.8
D+P	2.3
D+R+P	1.8
D+R+P+MFRM	1.2
D+R+P+MFRM+BCN	0.8

Table 2. Ablation study of refinement methods in MFRM.

Model	ACER(%)
D+R+P	1.8
D+R+P+Spatial Attention [53]	1.5
D+R+P+Channel Attention [54]	9.5
D+R+P+Non-local Attention [55]	12.7
D+R+P+Context-aware Reassembling	1.2

(ACER=2.1%) is obtained when $\sigma_r^2 = 1.0$ in BCN. We use this setting for the following experiments. According to the quantitative index shown in Fig. 6(b), BCN (ACER=2.1%) could decrease the ACER by half when compared with the results of ‘RGB.Baseline’ (ACER=4.1%). In order to validate whether the improvement from BCN is due to the extra parameters from ‘BilateralConvBlock’ in Fig. 4, we remove the bilateral operator in these blocks. However, as shown in Fig. 6(b) ‘BCN_w/o Bilateral’, it even works worse than the baseline (4.3% versus 4.1% ACER), indicating that the network easily overfits by introducing extra parameters without bilateral operators. We are also curious about the efficacy of bilateral residual term. After removing bilateral residual structure from BCN, the ACER index sharply changes from 2.1% (see ‘BCN’) to 4.0% (see ‘BCN_w/o Residual’). It implies that the micro-patterns from bilateral residual branch are also important for FAS task.

Influence of Various Input Types. As discussed in Sec. 3.1, the first solution is to adopt bilateral filtered images as network inputs. The results in Fig. 6(b) show that ‘Bi.Residual.Baseline’ with bilateral residual inputs performs better (0.4% ACER lower) than that of original RGB baseline. While combing the multi-level features from both original RGB and bilateral filtered images (i.e., ‘RGB+Bi_Base_Baseline’ and ‘RGB+Bi_Residual_Baseline’ in Fig. 6(b)), the performance further boosts. In contrast, BCN with only original RGB input outperforms the first solution methods for a large margin, implying that deep bilateral base and residual features in BCN are more robust.

Advantage of MFRM and Multi-head Supervision. Table 1 shows the ablation study about the network composition and supervision. ‘D’, ‘R’, ‘P’ are short for depth, reflection and patch heads, respectively. It is clear that

Table 3. The results of intra testing on four protocols of OULU-NPU.

Prot.	Method	APCER(%)↓	BPCER(%)↓	ACER(%)↓
1	GRADIANT [63]	1.3	12.5	6.9
	BASN [28]	1.5	5.8	3.6
	STASN [34]	1.2	2.5	1.9
	Auxiliary [14]	1.6	1.6	1.6
	FaceDs [11]	1.2	1.7	1.5
	FAS-TD [33]	2.5	0.0	1.3
	DeepPixBiS [26]	0.8	0.0	0.4
	Ours	0.0	1.6	0.8
2	DeepPixBiS [26]	11.4	0.6	6.0
	FaceDs [11]	4.2	4.4	4.3
	Auxiliary [14]	2.7	2.7	2.7
	BASN [28]	2.4	3.1	2.7
	GRADIANT [63]	3.1	1.9	2.5
	STASN [34]	4.2	0.3	2.2
	FAS-TD [33]	1.7	2.0	<u>1.9</u>
	Ours	2.6	0.8	1.7
3	DeepPixBiS [26]	11.7±19.6	10.6±14.1	11.1±9.4
	FAS-TD [33]	5.9±1.9	5.9±3.0	5.9±1.0
	GRADIANT [63]	2.6±3.9	5.0±5.3	3.8±2.4
	BASN [28]	1.8±1.1	3.6±3.5	2.7±1.6
	FaceDs [11]	4.0±1.8	3.8±1.2	3.6±1.6
	Auxiliary [14]	2.7±1.3	3.1±1.7	2.9±1.5
	STASN [34]	4.7±3.9	0.9±1.2	<u>2.8±1.6</u>
	Ours	2.8±2.4	2.3±2.8	2.5±1.1
4	DeepPixBiS [26]	36.7±29.7	13.3±14.1	25.0±12.7
	GRADIANT [63]	5.0±4.5	15.0±7.1	10.0±5.0
	Auxiliary [14]	9.3±5.6	10.4±6.0	9.5±6.0
	FAS-TD [33]	14.2±8.7	4.2±3.8	9.2±3.4
	STASN [34]	6.7±10.6	8.3±8.4	7.5±4.7
	FaceDs [11]	1.2±6.3	6.1±5.1	5.6±5.7
	BASN [28]	6.4±8.6	3.2±5.3	4.8±6.4
	Ours	2.9±4.0	7.5±6.9	5.2±3.7

multi-head supervision facilitates the network to learn more intrinsic features thus boost the performance. Furthermore, with both MFRM and multi-head supervision, our model is able to reduce ACER from baseline 4.1% to 1.2%. Ultimately, the full version of our method ‘D+R+P+MFRM+BCN’ achieves excellent performance with 0.8% ACER.

Impact of Refinement Methods in MFRM. We investigate four feature refinement methods in MFRM and the results are shown in Table 2. It is surprised that only spatial attention [53] (the second row) boosts the performance while SE block based channel attention [54] (the third row) and non-local block based self-attention[55] (the fourth row) perform poorly when domain shifts (e.g., illumination changes). We adopt context-aware reassembling as defaulted setting in MFRM as it can obtain more generalized features and improve baseline ‘D+R+P’ by 0.6% ACER. In summary, we use ‘D+R+P+MFRM+BCN (with $\sigma_r^2 = 1.0$)’ for all the following tests.

4.4 Intra Testing

The intra testing is carried out on both the OULU-NPU and the SiW datasets. We strictly follow the four protocols on OULU-NPU and three protocols on SiW for the evaluation. All compared methods including STASN [34] are trained without extra datasets for a fair comparison.

Results on OULU-NPU. As shown in Table 3, our proposed method ranks first or second on all the four protocols (0.4%, 1.7%, 2.5% and 5.2% ACER,

Table 4. The evaluation and comparison of the cross-type testing on SiW-M [29].

Method	Metrics(%)	Replay	Print	Mask Attacks				Makeup Attacks			Partial Attacks		Average		
				Half	Silicone	Trans.	Paper	Manne.	Obfusc.	Imperson.	Cosmetic	Funny Eye		Paper Glasses	Partial Paper
SVM+LBP [58]	APCER	19.1	15.4	40.8	20.3	70.3	0.0	4.6	96.9	35.3	11.3	53.3	58.5	0.6	32.8±29.8
	BPCER	22.1	21.5	21.9	21.4	20.7	23.1	22.9	21.7	12.5	22.2	18.4	20.0	22.9	21.0±2.9
	ACER	20.6	18.4	31.3	21.4	45.5	11.6	13.8	59.3	23.9	16.7	35.9	39.2	11.7	26.9±14.5
	EER	20.8	18.6	36.3	21.4	37.2	7.5	14.1	51.2	19.8	16.1	34.4	33.0	7.9	24.5±12.9
Auxiliary [14]	APCER	23.7	7.3	27.7	18.2	97.8	8.3	16.2	100.0	18.0	16.3	91.8	72.2	0.4	38.3±37.4
	BPCER	10.1	6.5	10.9	11.6	6.2	7.8	9.3	11.6	9.3	7.1	6.2	8.8	10.3	8.9± 2.0
	ACER	16.8	6.9	19.3	14.9	52.1	8.0	12.8	55.8	13.7	11.7	49.0	40.5	5.3	23.6±18.5
	EER	14.0	4.3	11.6	12.4	24.6	7.8	10.0	72.3	10.1	9.4	21.4	18.6	4.0	17.0±17.7
DTN [29]	APCER	1.0	0.0	0.7	24.5	58.6	0.5	3.8	73.2	13.2	12.4	17.0	17.0	0.2	17.1±23.3
	BPCER	18.6	11.9	29.3	12.8	13.4	8.5	23.0	11.5	9.6	16.0	21.5	22.6	16.8	16.6 ±6.2
	ACER	9.8	6.0	15.0	18.7	36.0	4.5	7.7	48.1	11.4	14.2	19.3	19.8	8.5	16.8 ±11.1
	EER	10.0	2.1	14.4	18.6	26.5	5.7	9.6	50.2	10.1	13.2	19.8	20.5	8.8	16.1± 12.2
Ours	APCER	12.4	5.2	8.3	9.7	13.6	0.0	2.5	30.4	0.0	12.0	22.6	15.9	1.2	10.3±9.1
	BPCER	13.2	6.2	13.1	10.8	16.3	3.9	2.3	34.1	1.6	13.9	23.2	17.1	2.3	12.2±9.4
	ACER	12.8	5.7	10.7	10.3	14.9	1.9	2.4	32.3	0.8	12.9	22.9	16.5	1.7	11.2±9.2
	EER	13.4	5.2	8.3	9.7	13.6	5.8	2.5	33.8	0.0	14.0	23.3	16.6	1.2	11.3±9.5

respectively), which indicates the proposed method performs well at the generalization of the external environment, attack mediums and input camera variation. Note that DeepPixBis [26] utilizes patch map for supervision but performing poorly in unseen attack mediums and camera types while BASN [26] exploits depth and reflection map as guidance but ineffective in unseen external environment. Our method works well for all 4 protocols as the extracted material-based features are intrinsic and generalized.

Results on SiW. We also compare our method with four state-of-the-art methods [14,34,33,28] on SiW dataset. Our method performs the best for all three protocols (0.36%, 0.11%, 2.45% ACER, respectively), revealing the excellent generalization capacity for 1) variations of face pose and expression, 2) variations of different spoof mediums, and 3) unknown presentation attack. More detailed results of each protocol are shown in *Appendix C*.

4.5 Inter Testing

To further validate whether our model is able to learn intrinsic features, we conduct cross-type and cross-dataset testing to verify the generalization capacity to unknown presentation attacks and unseen environment, respectively.

Cross-type Testing. Here we use CASIA-MFSD [59], Replay-Attack [60] and MSU-MFSD [61] to perform intra-dataset cross-type testing between replay and print attacks. Our proposed method achieves the best overall performance (96.77% AUC) among state-of-the-art methods [64,65,29,8], indicating the learned features generalized well among unknown attacks. More details can be found in *Appendix D*. Moreover, we also conduct cross-type testing on the latest SiW-M [29] dataset. As illustrated in Table 4, the proposed method achieves the best average ACER (11.2%) and EER (11.3%) among 13 attacks, which indicates our method actually learns material-based intrinsic patterns from rich kinds of material hence generalized well in unseen material type.

Cross-dataset Testing. In this experiment, we first train on the CASIA-MFSD and test on Replay-Attack, which is named as protocol CR. And then exchanging the training dataset and the testing dataset reciprocally, named protocol RC. As shown in Table 5, our proposed method has 16.6% HTER on protocol CR, outperforming all prior state-of-the-arts. For protocol RC, we also

Table 5. The results of cross-dataset testing between CASIA-MFSD and Replay-Attack. The evaluation metric is HTER(%).

Method	Protocol CR		Protocol RC	
	Train	Test	Train	Test
	CASIA-MFSD	Replay-Attack	Replay-Attack	CASIA-MFSD
LBP-TOP [66]	49.7			60.6
STASN [34]	31.5			30.9
Color Texture [3]	30.3			37.7
FaceDs [11]	28.5			41.1
Auxiliary [14]	27.6			28.4
BASN [28]	23.6			29.9
FAS-TD [33]	17.5			24.0
Ours	16.6			36.4

Table 6. Fine-grained material recognition in SiW-M dataset. The evaluation metric is accuracy (%).

Method	Live	Replay	Print	Mask	Makeup	Overall
ResNet50 (pre-trained) [51]	88.4	93.9	84.2	92.6	98.6	91.3
Patch (Ours)	91.5	98.0	87.7	95.9	73.9	90.1
Patch+BCN (Ours)	83.7	100.0	93.0	96.0	94.2	92.0
Patch+BCN+MFRM (Ours)	96.1	100.0	93.0	95.1	82.6	93.7

achieve comparable performance with 36.4% HTER. As our method is frame-level based, the performance might be further improved via introducing the temporal dynamic features in FAS-TD [33].

4.6 Analysis and Visualization.

In this subsection, we firstly conduct fine-grained material recognition on FAS dataset to prove that the learned features are material-based and intrinsic. And then we visualize and analyze the learned features.

Fine-grained Material Recognition. Face anti-spoofing is usually regarded as a binary classification problem despite we treat it as a binary material perception problem (i.e., structural facial skin versus others) in this paper. It is curious whether the model learns material-based intrinsic features despite it achieves state-of-the-art performance in most FAS datasets. As there are rich spoofing material types in SiW-M dataset, we separate it into five material categories, i.e., live, replay, print, mask and makeup, which are made of structural facial skin, plain glass, wrapped paper, structural fiber and foundation, respectively. Half samples of the each category are used for training and the remaining parts are utilized for testing. Only patch map supervision with five categories is utilized as the baseline because depth and reflection maps are not suitable for multi-category classification. As shown in Table 6, with the BCN and MFRM, the overall accuracy boosts by 1.9% and extra 1.7% respectively, outperforming ResNet50 [51] pre-trained in ImageNet. It implies intrinsic patterns among materials might be captured by BCN and MFRM.

Features Visualization. The low-level features of BCN and the predicted maps are visualized in Fig. 7. It is clear that the bilateral base and residual features between live and spoofing faces are quite different. For the bilateral base features of print attacks (see the 2nd and 3rd column in Fig. 7), the random noises in the hair region are obvious which is caused by the rough surface of the

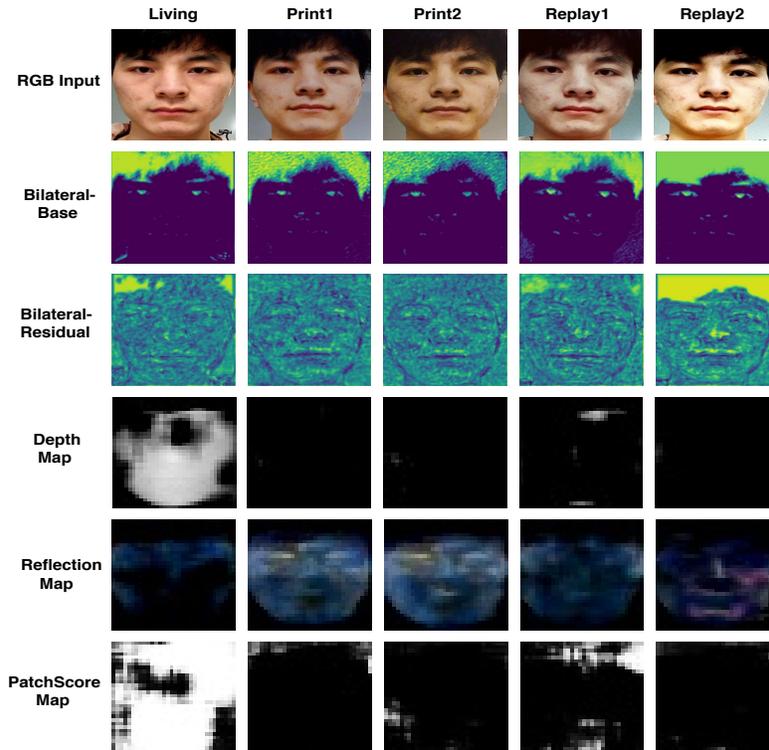


Fig. 7. Features visualization on live face (the first column) and spoofing faces (four columns to the right). The six rows represent the RGB images, low-level bilateral base features, low-level bilateral residual features in BCN, predicted depth maps, reflection maps and patch maps, respectively. Best view when zoom in.

paper material. Moreover, the micro- patterns in bilateral residual features reveal more details of facial outline in spoofing attacks but blurriness in live faces.

5 Conclusions

In this paper, we rephrase face anti-spoofing (FAS) task as a material recognition problem and combine FAS with classical human material perception [1]. To this end, Bilateral Convolutional Networks are proposed for capturing material-based bilateral macro- and micro- features. Extensive experiments are performed to verify the effectiveness of the proposed method. Our future works include: 1) to learn intrinsic material features via disentangling them with material-unrelated features (e.g., face id and face attribute features); and 2) to establish a more suitable cross-material based FAS benchmark.

Acknowledgment This work was supported by the Academy of Finland for project MiGA (grant 316765), ICT 2023 project (grant 328115), and Infotech Oulu. As well, the authors wish to acknowledge CSC IT Center for Science, Finland, for computational resources.

6 Appendix

A. Impact of Spatial Neighborhood Distance in DBO

The full version of deep bilateral operator (DBO) with spatial neighborhood distance term can be formulated as

$$DBO_{full}(\mathcal{F})_p = \frac{1}{k} \sum_{q \in \mathcal{F}} g_{\sigma_s}(\|p - q\|) g_{\sigma_r}(|\mathcal{F}_p - \mathcal{F}_q|) \mathcal{F}_q,$$

$$\text{with: } k = \sum_{q \in \mathcal{F}} g_{\sigma_s}(\|p - q\|) g_{\sigma_r}(|\mathcal{F}_p - \mathcal{F}_q|).$$

In this ablation study, the impact of spatial neighborhood distance $g_{\sigma_s}(\|p - q\|)$ would be evaluated. Here the default setting $\sigma_r^2 = 1.0$ is utilized. It can be seen from Fig. 8 that there are no improvements (2.4% and 2.1% ACER for with and without spatial neighborhood distance, respectively) when introducing spatial neighborhood distance term into BCN.

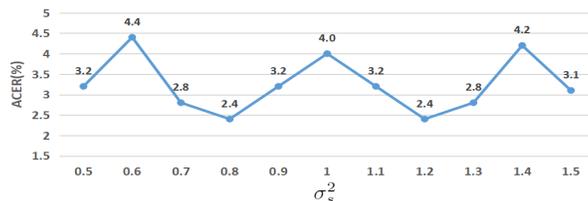


Fig. 8. Impact of σ_s in BCN.

Note that the ablation study about distance term σ_s here is not enough thus it might be a sub-optimal solution. The optimal hyperparameter setting could be found via strict grid search, which is one of our future works. Long-range spatial impact of distance term σ_s under large kernel size (e.g., 5x5 and 7x7) is also worth exploring in future.

B. Network Details of Multi-head Supervision

The detailed convolutional layers are illustrated in Fig. 9. With the supervision from three kinds of cues, the backbone network is able to learn more holistic material-based features.

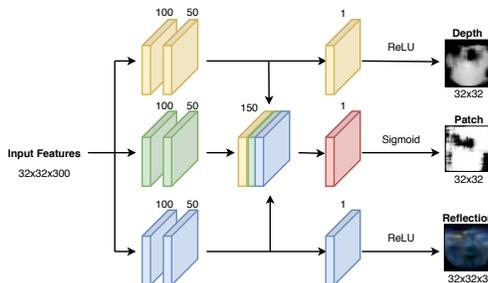


Fig. 9. Network structure of multi-head supervision. The number of filters are shown on top of each convolution, the size of all filters is 3×3 with stride 1.

C. Intra Testing Results on SiW

As shown in Table 7, the proposed method performs the best for all three protocols, revealing the excellent generalization capacity.

Table 7. The results of intra testing on three protocols of SiW [14].

Prot.	Method	APCER(%)	BPCER(%)	ACER(%)
1	Auxiliary [14]	3.58	3.58	3.58
	STASN [34]	–	–	1.00
	FAS-TD [33]	0.96	0.50	0.73
	BASN [28]	–	–	0.37
	Ours	0.55	0.17	0.36
2	Auxiliary [14]	0.57±0.69	0.57±0.69	0.57±0.69
	STASN [34]	–	–	0.28±0.05
	FAS-TD [33]	0.08±0.14	0.21±0.14	0.15±0.14
	BASN [28]	–	–	0.12±0.03
	Ours	0.08±0.17	0.15±0.00	0.11±0.08
3	STASN [34]	–	–	12.10±1.50
	Auxiliary [14]	8.31±3.81	8.31±3.80	8.31±3.81
	BASN [28]	–	–	6.45±1.80
	FAS-TD [33]	3.10±0.81	3.09±0.81	3.10±0.81
	Ours	2.55±0.89	2.34±0.47	2.45±0.68

D. Cross-type Testing on CASIA-MFSD, Replay-Attack and MSU-MFSD

In these cross-type testing, three datasets CASIA-MFSD [59], Replay-Attack [60] and MSU-MFSD [61] are utilized to perform intra-dataset cross-type testing between replay and print attacks. For instance, the second column ‘Video’ in Table 8 means that model should be trained from ‘Cut Photo’ and ‘Wrapped Photo’ while tested on ‘Video’. Table 8 shows that our proposed method achieves the best overall performance (96.77% AUC), indicating the learned features generalized well among unknown attacks.

Table 8. Cross-type testing on CASIA-MFSD, Replay-Attack, and MSU-MFSD. The evaluation metric is AUC (%).

Method	CASIA-MFSD [59]			Replay-Attack [60]			MSU-MFSD [61]			Overall
	Video	Cut Photo	Wrapped Photo	Video	Digital Photo	Printed Photo	Printed Photo	HR Video	Mobile Video	
OC-SVM+BBSIF [64]	70.74	60.73	95.90	84.03	88.14	73.66	64.81	87.44	74.69	78.68±11.74
SVM+LBP [58]	91.94	91.70	84.47	99.08	98.17	87.28	47.68	99.50	97.61	88.55±16.25
NN+LBP [65]	94.16	88.39	79.85	99.75	95.17	78.86	50.57	99.93	93.54	86.69±16.25
DTN [29]	90.0	97.3	97.5	99.9	99.9	99.6	81.6	99.9	97.5	95.9±6.2
AIM-FAS [8]	93.6	99.7	99.1	99.8	99.9	99.8	76.3	99.9	99.1	96.4±7.8
Ours	99.62	100.00	100.00	99.99	99.74	99.91	71.64	100.00	99.99	96.77±9.99

References

1. Sharan, L., Liu, C., Rosenholtz, R., Adelson, E.H.: Recognizing materials using perceptually inspired features. *International journal of computer vision* **103**(3) (2013) 348–371 [1](#), [3](#), [4](#), [5](#), [14](#)
2. Boulkenafet, Z., Komulainen, J., Hadid, A.: Face anti-spoofing based on color texture analysis. In: *IEEE international conference on image processing (ICIP)*. (2015) 2636–2640 [1](#), [3](#)
3. Boulkenafet, Z., Komulainen, J., Hadid, A.: Face spoofing detection using colour texture analysis. *IEEE Transactions on Information Forensics and Security* **11**(8) (2016) 1818–1830 [1](#), [13](#)
4. de Freitas Pereira, T., Anjos, A., De Martino, J.M., Marcel, S.: Lbp- top based countermeasure against face spoofing attacks. In: *Asian Conference on Computer Vision*. (2012) 121–132 [1](#), [3](#)
5. Komulainen, J., Hadid, A., Pietikainen, M.: Context based face anti-spoofing. In: *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*. (2013) 1–8 [1](#), [3](#)
6. Peixoto, B., Michelassi, C., Rocha, A.: Face liveness detection under bad illumination conditions. In: *ICIP, IEEE* (2011) 3557–3560 [1](#), [3](#)
7. Patel, K., Han, H., Jain, A.K.: Secure face unlock: Spoof detection on smartphones. *IEEE transactions on information forensics and security* **11**(10) (2016) 2268–2283 [1](#), [3](#)
8. Qin, Y., Zhao, C., Zhu, X., Wang, Z., Yu, Z., Fu, T., Zhou, F., Shi, J., Lei, Z.: Learning meta model for zero-and few-shot face anti-spoofing. *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)* (2020) [1](#), [3](#), [12](#), [16](#)
9. Wang, Z., Yu, Z., Zhao, C., Zhu, X., Qin, Y., Zhou, Q., Zhou, F., Lei, Z.: Deep spatial gradient and temporal depth learning for face anti-spoofing. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. (2020) 5042–5051 [1](#), [3](#)
10. Yu, Z., Zhao, C., Wang, Z., Qin, Y., Su, Z., Li, X., Zhou, F., Zhao, G.: Searching central difference convolutional networks for face anti-spoofing. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. (2020) 5295–5305 [1](#), [3](#)
11. Jourabloo, A., Liu, Y., Liu, X.: Face de-spoofing: Anti-spoofing via noise modeling. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. (2018) 290–306 [1](#), [3](#), [11](#), [13](#)
12. Lin, B., Li, X., Yu, Z., Zhao, G.: Face liveness detection by rppg features and contextual patch-based cnn. In: *Proceedings of the 2019 3rd International Conference on Biometric Engineering and Applications, ACM* (2019) 61–68 [2](#), [3](#)
13. Li, X., Komulainen, J., Zhao, G., Yuen, P.C., Pietikäinen, M.: Generalized face anti-spoofing by detecting pulse from face videos. In: *2016 23rd International Conference on Pattern Recognition (ICPR), IEEE* (2016) 4244–4249 [2](#), [3](#)
14. Liu, Y., Jourabloo, A., Liu, X.: Learning deep models for face anti-spoofing: Binary or auxiliary supervision. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2018) 389–398 [2](#), [3](#), [5](#), [6](#), [8](#), [9](#), [10](#), [11](#), [12](#), [13](#), [16](#)
15. Liu, S.Q., Lan, X., Yuen, P.C.: Remote photoplethysmography correspondence feature for 3d mask face presentation attack detection. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. (2018) 558–573 [2](#)
16. Atoum, Y., Liu, Y., Jourabloo, A., Liu, X.: Face anti-spoofing using patch and depth-based cnns. In: *2017 IEEE International Joint Conference on Biometrics (IJCB)*. (2017) 319–328 [2](#), [3](#)

17. Tan, X., Li, Y., Liu, J., Jiang, L.: Face liveness detection from a single image with sparse low rank bilinear discriminative model. In: European Conference on Computer Vision, Springer (2010) 504–517 [2](#), [4](#)
18. Li, L., Xia, Z., Jiang, X., Ma, Y., Roli, F., Feng, X.: 3d face mask presentation attack detection based on intrinsic image analysis. arXiv preprint arXiv:1903.11303 (2019) [2](#), [4](#)
19. Boulkenafet, Z., Komulainen, J., Hadid, A.: Face antispoofing using speeded-up robust features and fisher vector encoding. *IEEE Signal Processing Letters* **24**(2) (2017) 141–145 [3](#)
20. Komulainen, J., Hadid, A., Pietikäinen, M.: Face spoofing detection using dynamic texture. In: Asian Conference on Computer Vision, Springer (2012) 146–157 [3](#)
21. Siddiqui, T.A., Bharadwaj, S., Dhamecha, T.I., Agarwal, A., Vatsa, M., Singh, R., Ratha, N.: Face anti-spoofing with multifeature videolet aggregation. In: 2016 23rd International Conference on Pattern Recognition (ICPR), IEEE (2016) 1035–1040 [3](#)
22. Pan, G., Sun, L., Wu, Z., Lao, S.: Eyeblick-based anti-spoofing in face recognition from a generic webcam. In: IEEE International Conference on Computer Vision. (2007) 1–8 [3](#)
23. Yu, Z., Qin, Y., Xu, X., Zhao, C., Wang, Z., Lei, Z., Zhao, G.: Auto-fas: Searching lightweight networks for face anti-spoofing. In: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE (2020) 996–1000 [3](#)
24. Li, L., Feng, X., Boulkenafet, Z., Xia, Z., Li, M., Hadid, A.: An original face anti-spoofing approach using partial convolutional neural network. In: IPTA. (2016) 1–6 [3](#)
25. Patel, K., Han, H., Jain, A.K.: Cross-database face antispoofing with robust feature representation. In: Chinese Conference on Biometric Recognition. (2016) 611–619 [3](#)
26. George, A., Marcel, S.: Deep pixel-wise binary supervision for face presentation attack detection. In: International Conference on Biometrics. Number CONF (2019) [3](#), [11](#), [12](#)
27. Yu, Z., Qin, Y., Li, X., Wang, Z., Zhao, C., Lei, Z., Zhao, G.: Multi-modal face anti-spoofing based on central difference networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. (2020) 650–651 [3](#)
28. Kim, T., Kim, Y., Kim, I., Kim, D.: Basn: Enriching feature representation using bipartite auxiliary supervisions for face anti-spoofing. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. (2019) 0–0 [3](#), [8](#), [11](#), [12](#), [13](#), [16](#)
29. Liu, Y., Stehouwer, J., Jourabloo, A., Liu, X.: Deep tree learning for zero-shot face anti-spoofing. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2019) 4680–4689 [3](#), [9](#), [12](#), [16](#)
30. Jia, Y., Zhang, J., Shan, S., Chen, X.: Single-side domain generalization for face anti-spoofing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2020) 8484–8493 [3](#)
31. Wang, G., Han, H., Shan, S., Chen, X.: Cross-domain face presentation attack detection via multi-domain disentangled representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2020) 6678–6687 [3](#)

32. Shao, R., Lan, X., Li, J., Yuen, P.C.: Multi-adversarial discriminative deep domain generalization for face presentation attack detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2019) 10023–10031 [3](#)
33. Wang, Z., Zhao, C., Qin, Y., Zhou, Q., Lei, Z.: Exploiting temporal and depth information for multi-frame face anti-spoofing. arXiv preprint arXiv:1811.05118 (2018) [3](#), [8](#), [11](#), [12](#), [13](#), [16](#)
34. Yang, X., Luo, W., Bao, L., Gao, Y., Gong, D., Zheng, S., Li, Z., Liu, W.: Face anti-spoofing: Model matters, so does data. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2019) [3](#), [11](#), [12](#), [13](#), [16](#)
35. Lin, C., Liao, Z., Zhou, P., Hu, J., Ni, B.: Live face verification with multiple instantiated local homographic parameterization. In: IJCAI. (2018) 814–820 [3](#)
36. Yu, Z., Li, X., Zhao, G.: Remote photoplethysmograph signal measurement from facial videos using spatio-temporal networks. arXiv preprint arXiv:1905.02419 (2019) [3](#)
37. Yu, Z., Li, X., Niu, X., Shi, J., Zhao, G.: Autohr: A strong end-to-end baseline for remote heart rate measurement with neural searching. arXiv preprint arXiv:2004.12292 (2020) [3](#)
38. Yu, Z., Peng, W., Li, X., Hong, X., Zhao, G.: Remote heart rate measurement from highly compressed facial videos: an end-to-end deep learning solution with video enhancement. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 151–160 [3](#)
39. Maloney, L.T., Brainard, D.H.: Color and material perception: Achievements and challenges. *Journal of Vision* **10**(9) (2010) 19–19 [4](#)
40. Fleming, R.W.: Visual perception of materials and their properties. *Vision research* **94** (2014) 62–75 [4](#)
41. Nishida, S.: Image statistics for material perception. *Current Opinion in Behavioral Sciences* **30** (2019) 94–99 [4](#)
42. Adelson, E.H.: On seeing stuff: the perception of materials by humans and machines. In: Human vision and electronic imaging VI. Volume 4299., International Society for Optics and Photonics (2001) 1–12 [4](#)
43. Varma, M., Zisserman, A.: A statistical approach to material classification using image patch exemplars. *IEEE transactions on pattern analysis and machine intelligence* **31**(11) (2008) 2032–2047 [4](#)
44. Jiang, X., Du, J., Sun, B., Feng, X.: Deep dilated convolutional network for material recognition. In: 2018 Eighth International Conference on Image Processing Theory, Tools and Applications (IPTA), IEEE (2018) 1–6 [4](#)
45. Ling, S., Callet, P.L., Yu, Z.: The role of structure and textural information in image utility and quality assessment tasks. *Electronic Imaging* **2018**(14) (2018) 1–13 [4](#)
46. Deng, B.W., Yu, Z.T., Ling, B.W., Yang, Z.: Video quality assessment based on features for semantic task and human material perception. In: 2016 IEEE International Conference on Consumer Electronics-China (ICCE-China), IEEE (2016) 1–4 [4](#)
47. Li, L., Xia, Z., Jiang, X., Ma, Y., Roli, F., Feng, X.: 3d face mask presentation attack detection based on intrinsic image analysis. *IET Biometrics* **9**(3) (2020) 100–108 [4](#)
48. Tomasi, C., Manduchi, R.: Bilateral filtering for gray and color images. In: *Iccv*. Volume 98. (1998) [2](#) [4](#)
49. Durand, F., Dorsey, J.: Fast bilateral filtering for the display of high-dynamic-range images. In: *ACM transactions on graphics (TOG)*. Volume 21., ACM (2002) 257–266 [5](#)

50. Paris, S., Durand, F.: A fast approximation of the bilateral filter using a signal processing approach. In: European conference on computer vision, Springer (2006) 568–580 [5](#)
51. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2016) 770–778 [6](#), [13](#)
52. Wang, J., Chen, K., Xu, R., Liu, Z., Loy, C.C., Lin, D.: Carafe: Content-aware reassembly of features. arXiv preprint arXiv:1905.02188 (2019) [6](#)
53. Woo, S., Park, J., Lee, J.Y., So Kweon, I.: Cbam: Convolutional block attention module. In: Proceedings of the European Conference on Computer Vision (ECCV). (2018) 3–19 [7](#), [10](#), [11](#)
54. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2018) 7132–7141 [7](#), [10](#), [11](#)
55. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 7794–7803 [7](#), [10](#), [11](#)
56. Feng, Y., Wu, F., Shao, X., Wang, Y., Zhou, X.: Joint 3d face reconstruction and dense alignment with position map regression network. In: Proceedings of the European Conference on Computer Vision (ECCV). (2017) [8](#)
57. Zhang, X., Ng, R., Chen, Q.: Single image reflection separation with perceptual losses. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 4786–4794 [8](#)
58. Boulkenafet, Z., Komulainen, J., Li, L., Feng, X., Hadid, A.: Oulu-npu: A mobile face presentation attack database with real-world variations. In: FGR. (2017) 612–618 [9](#), [12](#), [16](#)
59. Zhang, Z., Yan, J., Liu, S., Lei, Z., Yi, D., Li, S.Z.: A face antispoofing database with diverse attacks. In: ICB. (2012) 26–31 [9](#), [12](#), [16](#)
60. Chingovska, I., Anjos, A., Marcel, S.: On the effectiveness of local binary patterns in face anti-spoofing. In: Biometrics Special Interest Group. (2012) 1–7 [9](#), [12](#), [16](#)
61. Wen, D., Han, H., Jain, A.K.: Face spoof detection with image distortion analysis. IEEE Transactions on Information Forensics and Security **10**(4) (2015) 746–761 [9](#), [12](#), [16](#)
62. international organization for standardization: Iso/iec jtc 1/sc 37 biometrics: Information technology biometric presentation attack detection part 1: Framework. In: <https://www.iso.org/obp/ui/iso>. (2016) [9](#)
63. Boulkenafet, Z., Komulainen, J., Akhtar, Z., Benlamoudi, A., Samai, D., Bekhouche, S.E., Ouafi, A., Dornaika, F., Taleb-Ahmed, A., Qin, L., et al.: A competition on generalized software-based face presentation attack detection in mobile scenarios. In: 2017 IEEE International Joint Conference on Biometrics (IJCB), IEEE (2017) 688–696 [11](#)
64. Arashloo, S.R., Kittler, J., Christmas, W.: An anomaly detection approach to face spoofing detection: A new formulation and evaluation protocol. IEEE Access **5** (2017) 13868–13882 [12](#), [16](#)
65. Xiong, F., AbdAlmageed, W.: Unknown presentation attack detection with face rgb images. In: 2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS), IEEE (2018) 1–9 [12](#), [16](#)
66. de Freitas Pereira, T., Anjos, A., De Martino, J.M., Marcel, S.: Can face anti-spoofing countermeasures work in a real world scenario? In: 2013 international conference on biometrics (ICB). (2013) 1–8 [13](#)