

Uncertainty-Aware Physically-Guided Proxy Tasks for Unseen Domain Face Anti-spoofing

Junru Wu¹, Xiang Yu², Buyu Liu², Zhangyang Wang³, Manmohan Chandraker²
¹Texas A&M University, ²NEC Laboratories America, ³University of Texas at Austin
 sandboxmaster@tamu.edu, {xiangyu, buyu, manu}@nec-labs.com, atlaswang@utexas.edu

Abstract

Face anti-spoofing (FAS) seeks to discriminate genuine faces from fake ones arising from any type of spoofing attack. Due to the wide varieties of attacks, it is implausible to obtain training data that spans all attack types. We propose to leverage physical cues to attain better generalization on unseen domains. As a specific demonstration, we use physically guided proxy cues such as depth, reflection, and material to complement our main anti-spoofing (a.k.a liveness detection) task, with the intuition that genuine faces across domains have consistent face-like geometry, minimal reflection, and skin material. We introduce a novel uncertainty-aware attention scheme that independently learns to weigh the relative contributions of the main and proxy tasks, preventing the over-confident issue with traditional attention modules. Further, we propose attribute-assisted hard negative mining to disentangle liveness-irrelevant features with liveness features during learning. We evaluate extensively on public benchmarks with intra-dataset and inter-dataset protocols. Our method achieves the superior performance especially in unseen domain generalization for FAS.

1. Introduction

With growing prevalence of face recognition, it is increasingly subject to a wide variety of spoofing attacks. Thus, face anti-spoofing or liveness detection is emerging as an essential precursor, with the key need being robust to various attacks that are possibly unseen previously and drastically different in appearance from training data. This is an extremely challenging problem due to the fact that sophisticated spoofs might arise from similar camera and lighting setups as genuine inputs, leading to only subtle differences in appearance. On the other hand, types of attacks range from printed photos to facial masks, which makes it laborious to obtain exhaustive training data for anti-spoofing task.

In this paper, we aim to address the face anti-spoofing problem on *unseen* domains or attack types, where neither

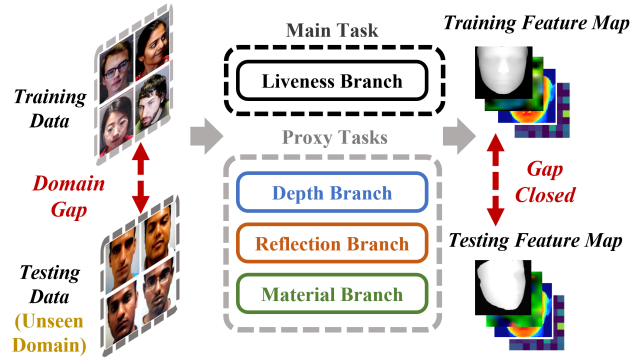


Figure 1: Traditional methods only focus on real/spoofing binary classification, which results in sensitive prediction. By introducing the physical cues of depth, material and reflection as proxy tasks, our method largely close the domain gap from the training data to the unseen testing data and thus boost the performance reliably.

definition of the attack types nor training data under supervised or unsupervised condition is available. To achieve this, we derive inspiration from physical cues that establish a commonality for genuine inputs and distinction from fake ones. We refer to the estimation of these cues as *proxy tasks*, performed in conjunction with the main task, a.k.a appearance-based liveness detection. While the proposed formulation is general and physical cues can be arbitrary, we focus on depth estimation, reflection detection and material classification as our proxy tasks. Intuitively, we expect genuine faces to constitute face-like geometry and present skin as the material, while several presentation attacks might violate at least one of those conditions. As a consequence, incorporating such proxy tasks enables to generalize the shared cues to unseen domains or attack types.

We bring the insights from single-image based face reconstruction using 3D morphable models (3DMM) [7] for depth proxy, single-image based material recognition trained on large-scale datasets [5] for material proxy, and a single image reflection separation model [58] to provide

the pseudo labels for the reflection proxy. A shared encoder is trained across the main and proxy tasks to transfer the insights into our deep appearance-based liveness detection problem. In contrast to existing work [34] that incorporates depth cue to regularize sensitive binary liveness task, our proxy tasks are more general in the sense of considering more physical cues such as material, to gear towards a physically meaningful way to handle unseen domains. Meanwhile, we organize the proxy tasks into a multi-channel learning framework to provide a more robust detection with an attention aggregation. Note that domain adaptation is not applicable in our setting, since we assume that even unlabeled training data is not available, which cannot define the target domain.

Besides proxy tasks, we also leverage a pretext task in the form of face recognition, which is usually regularized by large scale labeled datasets and expected to provide high-level shared face analysis feature representation for liveness detection. We thereby provide recipes for pre-training on face recognition that allows better generalization to unseen domains for the liveness task, as well as multi-channel training with liveness and proxy tasks. We conduct extensive experiments on five publicly available benchmarks. In each case, we demonstrate not only state-of-the-art results, but also that judicious use of pretext and proxy tasks allows better generalization of liveness detection to unseen domains. Besides, the multi-task learning can result in channel conflict as the liveness feature is ideally invariant to identity information where the pretext task and our liveness task share the same network. To this end, we leverage the attribute information in the proxy data to conduct a triplet metric learning based mining, expecting to better disentangle the non-liveness information from the learned feature and thus boosts the liveness detection.

To better exploit multiple physically meaningful resources, we further holistically weigh the relative contributions of the main task and various proxy tasks with an uncertainty-aware attention module. Traditional attention modules are jointly optimized with all tasks and might cause the notorious over-confident issue due to training data bias. While our uncertainty-aware attention is designed to independently estimate the tasks' variance, which does not capture feature fitness to the task but rather focusing on its deviation to the estimated mean or termed uncertainty of the feature estimation. This property ensures less bias in uncertainty-aware attention module thus captures the property of input images better.

In summary, we propose the following contributions:

- We propose three physical-cue guided proxy tasks including depth, material and reflection, which share the commonality across domains to enable the unseen domain anti-spoofing.
- We leverage an uncertainty-aware attention module

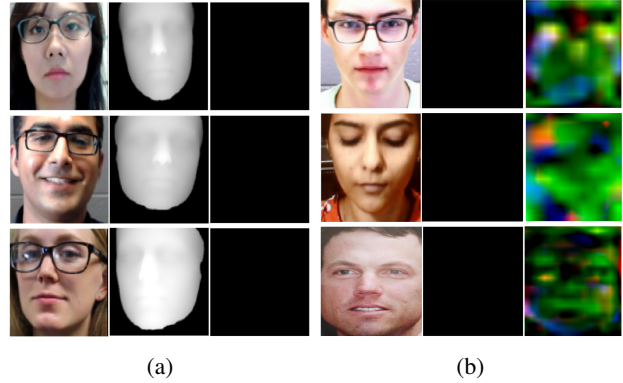


Figure 2: (a) From Left to Right: examples of face image, depth map, reflection map for genuine faces (b) From Left to Right: examples of face image, depth map, reflection map for spoof faces.

to effectively combine the main and proxy tasks and boost the performance.

- We propose an attribute-assisted mining scheme to make sure liveness-irrelevant features are properly disentangled and only liveness features are learned.
- We conduct an extensive evaluation with both intra-dataset and inter-dataset protocols including the latest attribute-rich CelebA-Spoof dataset, highlighting our framework's better performance in unseen domain generalization for FAS.

2. Related Work

We categorize face anti-spoofing literature into physical cue based, feature learning based methods, and whether they address the unseen spoofing attacks.

Physical Cue based Anti-spoofing: Early research on anti-spoofing leverages the physical cues, e.g. head movement [27] and eye blinking [42], to indicate the genuineness. These methods can be simply spoofed by printing faces with eye region cut, or wearing a facial mask and moving head. By analyzing the lighting cooperated from different reflection, the remote Photoplethysmography (r-PPG) [15, 8, 33, 41] is proposed to identify spoofing attacks with material information. However, this type of methods require the imaging quality to be high as the lighting measurement is less tolerated to noise. Combining with CNN, depth is proposed [34, 46] as an auxiliary task that enables less sensitive and more explainable training. [26] introduced reflection map as a supplement to depth for bipartite auxiliary supervision. This method is limited in generalizing to other spoofing resources since their depth, reflection and r-PPG are trained on a single dataset. Instead, our depth, reflection and material channels are guided by models trained on large scale datasets, i.e., a 3DMM based depth regression model [17], a single image reflection separation model [58] and a material classification model [5], which

Method	[35]	[44]	[34]	[46]	Ours
Physical Cues	×	Blink	Depth, rPPG	Depth	Depth, Reflection, Material
Temporal information	×	✓	✓	×	×
Multiple-domain data	✓	×	×	×	×
Cross-domain evaluation	✓	✓	✓	✓	✓

Table 1: Unseen domain anti-spoofing methods comparison with robust feature [44], anomaly detection [2], binary or auxiliary supervision [34], deep tree learning (DTL) [35] and MADDG [46]. “✓” means applying and “×” means not applying.

substantially improves the generalization.

Learning based Anti-spoofing: The handcrafted features, e.g., HoG [28, 23], SIFT [43] and LBP [9, 45, 39, 52, 57] are explored in early literature. Such binary classification achieves good performance but is restricted to some defined domains. Meanwhile, those methods do not consider environment variations, i.e., lighting, color tone or pose change. To this end, the HSV and YCbCr [9], Fourier transform [31] and image low-rank decomposition [51] methods are also explored. Some other works [1, 4, 49, 16, 61] utilize the temporal information assuming videos are available. Later, deep learning based features [56, 16, 32, 3, 22, 34, 44] are utilized [56, 16, 32, 3, 22, 34, 44] and achieve better performance. Notice that both [3] and [34] leverage texture and depth, which seem to be close to our setting. However, instead of directly exploring texture, we formulate the texture into a more physically consistent cue, the material, and set up the material classification task to avoid rPPG calculation. Moreover, our method is single image based which does not require the temporal information, thus reducing the run-time and model complexity.

Unseen Domain Anti-spoofing: Methods that explore handcrafted features can deal with unseen spoofing attacks as these features are independent from attack types. However, due to the limitation of feature representation power, they cannot generalize well. A method comparison is listed in Table 1. Patel et al. [44] propose to combine deep features with eye blinking cues for cross-dataset spoofing detection. There are works [2, 54] formulating the anti-spoofing task into an anomaly or outlier detection, which highly rely on the definition of genuine samples. To alleviate this, Liu et al. [35] propose a zero-shot learning solution, whereas the unseen attacks are assigned to the most similar attacks pre-defined in the database. These unseen attacks are wildly variant and leave the chance that the attacks are heavy outliers. Shao et al. [46] formulate a domain generalization approach to improve the generalization ability, which depends on the number and diversity of the seen domains, i.e., biased or long-tailed observed domains would degrade the performance. Different from [35, 46], we propose the physical cue based proxy tasks that are less dependent on the data distribution, which generally could be more stable and consistent across seen and unseen attacks.

Uncertainty Analysis: Uncertainty provides an effective measurement for model/data reliability [30, 50, 24, 38]. It

has been widely applied in many vision tasks such as classification [21], semantic segmentation [25] and face recognition [6]. Our method follows the setting of [25] by leveraging multiple tasks, but in a complete different problem as face anti-spoofing rather than segmentation. We consider our proxy tasks are orthogonal to each other, whereas in [25] those semantic tasks are strongly correlated. To the best of our knowledge, we are the first to leverage uncertainty in face anti-spoofing tasks.

3. Proposed Approach

In this section, we firstly introduce the shared feature extractor learning by incorporating the pretext task face recognition. Then, the physical cue based proxy tasks, i.e., depth estimation, reflection detection and material classification, are introduced as the spoofing attack detection anchors. Finally, an uncertainty-aware attention module is proposed to aggregate the proxy channels for optimal performance.

3.1. Shared Feature Representation Learning

As shown in Figure 3, our framework consists of multiple channels of pretext and proxy tasks. Separating each single task with independent CNNs results in network redundancy. Moreover, the separated CNNs cannot leverage the rich information from the other tasks, where hyper-column [18] and deeply-supervised net [29] have shown a highly integrated framework for multiple tasks is beneficial. To this end, we propose to use a single feature extractor Φ to provide the shared feature for all the downstream tasks.

The shared features should provide high-level task-specific yet general information for downstream tasks such that we neither drift away from original tasks nor learn only task-driven representations. Among the face analysis applications, face recognition is a promising pretext as it is usually trained with large-scale data including millions of identities, which guarantees the robustness as well as the discriminability. Other candidates such as facial attribute classification, expression recognition or spoofing detection are not general or robust, as each of the tasks conduct a 10-way or 2-way classification, which can be sensitive or easily overfitting [34]. Thus, to initialize the feature extractor Φ , we apply face recognition as our pretext task.

Denoting input image as \mathbf{x}_r , \mathbf{x}_v and \mathbf{x}_m for recognition, spoofing and material data respectively. After the shared feature extractor Φ , the pretext task applies a filter Ψ_r to refine the face identity feature. The loss is:

$$\mathcal{L}_r = - \sum_i \mathbb{1}(y_i) \log \frac{\exp(\mathbf{w}_i \Psi_r \circ \Phi(\mathbf{x}_r))}{\sum_j \exp(\mathbf{w}_j \Psi_r \circ \Phi(\mathbf{x}_r))} \quad (1)$$

where y_i is the ground truth label for identity i . j varies across the whole number of identities. \mathbf{w}_i is the i_{th} separation hyper-plane of the classifier. \circ denotes the sequential network flow.

\mathbf{x}_r - image (recognition)
 \mathbf{x}_v - image (spoofing)
 \mathbf{x}_m - image (material data)

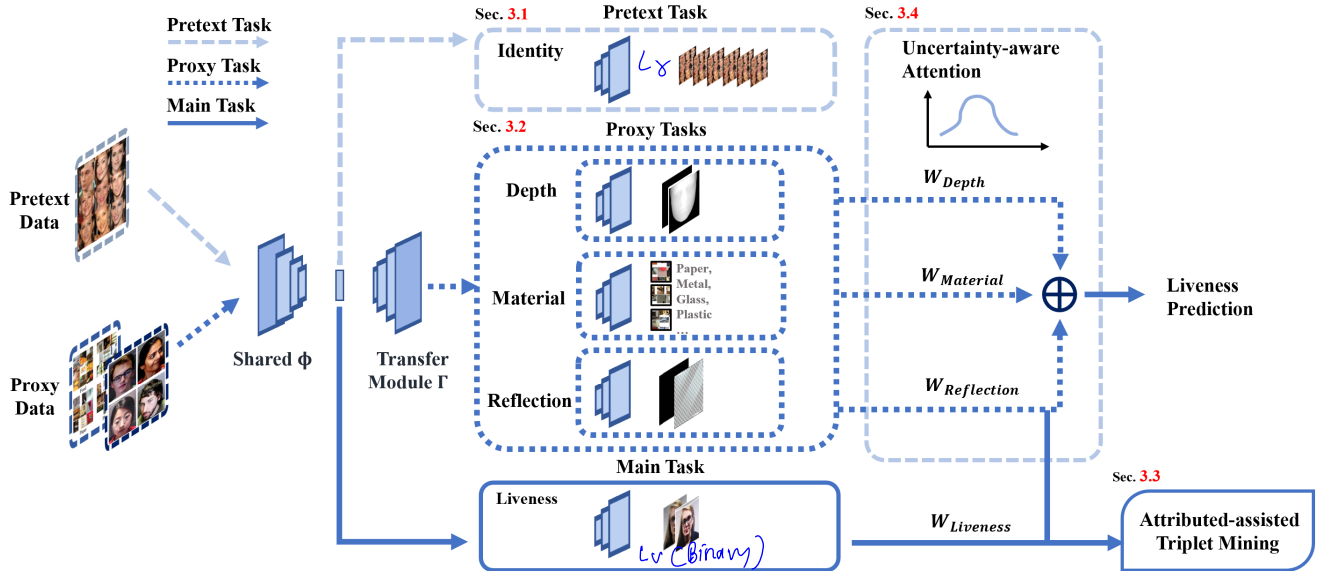


Figure 3: The proposed framework consists of the pretext task “face recognition” (Sec. 3.1), the proxy tasks (Sec. 3.2) “depth estimation”, “material prediction”, “reflection detection” and the main task “liveness detection”. A novel triplet mining regularization (Sec. 3.3) is proposed to better disentangle the liveness feature and an uncertainty-aware attention module (Sec. 3.4) aggregates the channel-wise results for boosted performance.

3.2. Multi-channel Proxy Task Learning

We introduce a transfer module Γ to adapt the rich feature extracted by Φ into the spoofing detection related tasks. Directly utilizing the shared feature leads to sub-optimal prediction as it incorporates unrelated face recognition cues, which may serve as noise. Similar to the pretext task, we set up multiple channels for our proxy tasks, i.e., liveness detection Ψ_v , depth estimation Ψ_d , reflection detection Ψ_r and material prediction Ψ_m .

Liveness Detection Main Task: The spoofing detection is a well-known binary classification task. The input is spoofing face \mathbf{x}_v . After shared extractor and feature transfer module, we set up the spoofing detection channel filter Ψ_v to conduct the binary classification task, in which we adopt binary cross entropy loss as the objective:

$$\mathcal{L}_v = -y_v \log(\mathbf{p}(\mathbf{z})) - (1 - y_v) \log(1 - \mathbf{p}(\mathbf{z})) \quad (2)$$

$$\mathbf{p}(\mathbf{z}) = \frac{\exp(\tilde{\mathbf{w}}_0 \mathbf{z})}{\exp(\tilde{\mathbf{w}}_0 \mathbf{z}) + \exp(\tilde{\mathbf{w}}_1 \mathbf{z})} \quad (3)$$

where \mathbf{y}_v is the ground truth of spoofing or genuine, $\mathbf{z} = \Psi_v \circ \Phi(\mathbf{x}_v)$ denotes the spoofing detection feature after the spoofing detection filter Ψ_v . $\tilde{\mathbf{w}}_0$ and $\tilde{\mathbf{w}}_1$ are the separation hyper-planes of the binary classifier. Likelihood of being spoofing sample $\mathbf{p}(\mathbf{z})$ is estimated via a softmax operation in Equation 3.

Depth Proxy Task: We believe the physical cues should share similar characteristics for genuine faces across different attack types or spoofing datasets. Thus the depth pre-

dition should also be consistent. We aim to predict the per-pixel depth map given the input face image. We leverage an hourglass network structure to conduct this regression problem, which has been proved effective in key point detection [40] and image segmentation [12]. To prepare the ground truth depth map \mathbf{d}_{GT} , we apply a 3D face shape reconstruction algorithm [17] offline to estimate the dense point cloud for the face images. As for genuine face image, we utilize the predicted depth as its ground truth depth map, where background is set as 0. For 2D spoofing face images, according to their attack types, i.e., display screen or paper, we know that the actual depth is from a flat plane of either screen or paper. Thus, we manually set the spoofing ground truth depth to be all 0. The absolute depth is unnecessary since we only focus on the relative face geometry. We show some examples of generated depth map results in Figure 2. During training, an l_1 -based reconstruction loss is applied:

$$\mathcal{L}_d = \|\Psi_d \circ \Gamma \circ \Phi(\mathbf{x}_v) - \mathbf{d}_{GT}\|_1 \quad (4)$$

where Ψ_d is the hourglass net depth estimation module, Γ is the feature transfer module and \mathbf{d}_{GT} is the ground truth depth. Notice that for depth estimation, we input the spoofing data \mathbf{x}_v with the augmented ground truth depth map. We do not utilize extra depth data for this channel.

Reflection Proxy Task: Reflection is another useful physical cue that indicates the genuine faces, as non-skin materials inevitably show abnormal reflection compared to skin. As a result, for spoofing face, we use a single image reflection separation model [58] to generate the reflection map,

while for genuine face, we set the it to zero denoting no reflection is present. Visual examples of generated reflection maps are in Figure 2. During training, we push the predicted reflection map to be close to the pseudo ground truth under l_1 -based constraint:

$$\mathcal{L}_r = \|\Psi_r \circ \Gamma \circ \Phi(\mathbf{x}_v) - \mathbf{R}_{GT}\|_1 \quad (5)$$

where Ψ_r is the hourglass net reflection estimation module and \mathbf{r}_{GT} is the ground truth reflection map.

Material Proxy Task: Though reflection in some way indicates the material information, we explicitly introduce material as another proxy task to leverage the correlation among the multiple tasks, expecting to benefit from the multi-task learning. The physical insight for material in liveness detection is that skin across different spoofing attacks or spoofing datasets should remain similar RGB information. We automatically obtain the material type for face spoofing data according to its attack type. For instance, we denote the material class of screen display and paper print as “glass” and “paper” respectively. In this way, we actually unify the material type towards the general material recognition [5].

Notice that the number of material types in spoofing data can be limited, which may encounter the same sensitivity issue as the binary spoofing detection task. To this end, we introduce the general material recognition data to anchor the material feature space from being collapsed. Specifically, the general material recognition and our spoofing data material recognition share all the network structures except the last classifier layer. As in general material recognition, there are 23 defined categories [5], such as brick, metal, plastic, skin, glass, etc. We set up a 23-way classifier \mathbf{C}_g for the general material recognition and a 3-way classifier \mathbf{C}_v for our spoofing data material recognition. A multi-source scheme is proposed to train the modules of Φ , Γ and Ψ_m jointly. Denoting the feature $\mathbf{f} = \Psi_m \circ \Gamma \circ \Phi(\mathbf{x})$, a combined multi-class softmax loss is applied to train \mathbf{C}_g and \mathbf{C}_v :

$$\mathcal{L}_m = - \sum_{i=1}^{23} \mathbb{1}(l_i) \log \frac{\exp(\omega_i \Psi_m \circ \Gamma \circ \Phi(\mathbf{x}_m))}{\sum_j \exp(\omega_j \Psi_m \circ \Gamma \circ \Phi(\mathbf{x}_m))} - \sum_{i=1}^3 \mathbb{1}(l_i) \log \frac{\exp(\tilde{\omega}_i \Psi_m \circ \Gamma \circ \Phi(\mathbf{x}_v))}{\sum_j \exp(\tilde{\omega}_j \Psi_m \circ \Gamma \circ \Phi(\mathbf{x}_v))} \quad (6)$$

where l_i is the material ground truth label, ω_i and ω_j , $\tilde{\omega}_i$ and $\tilde{\omega}_j$ are the separation hyper-planes for \mathbf{C}_g and \mathbf{C}_v respectively. By alternatively feeding the material and spoofing data, we guarantee that \mathbf{f} is generalized for not only the standard material recognition, but also the material recognition for face spoofing data.

3.3. Attributed-assisted Triplet Mining

To better disentangle the liveness feature apart from identity information and other facial attributes information, we leverage the metric learning to regularize the feature representation learning. Specifically, given the input \mathbf{x}_v^i , we would expect the following loss to be minimized such that the identity information can be decoupled from the liveness feature.

$$\mathcal{L}_{tid} = [\|\Phi(x_v^{i,j}) - \Phi(x_v^{i,k})\|^2 - \|\Phi(x_v^{i,j}) - \Phi(x_v^{h,k})\|^2 + m_1]_+ \quad (7)$$

Anchor *Anchor encodey* → *positive sample encodey* } *Triplet loss*

$x_v^{i,j}$ means the j^{th} liveness sample from identity i , while x_v^h simply means the liveness sample from other identities as a negative sample.

Similarly for other face attributes introduced in CelebA [36, 59], we believe the orthogonality can be preserved if those attribute information is disentangled from the liveness information.

$$\mathcal{L}_{ta} = [\|\Phi(x_v^{a_i,j}) - \Phi(x_v^{a_i,k})\|^2 - \|\Phi(x_v^{a_i,j}) - \Phi(x_v^{a_h,k})\|^2 + m_2]_+ \quad (8)$$

a_i indicates an attribute label and a_h indicates a different attribute label for the negative sample. m_1 and m_2 here are the margin hyper-parameter set to squeeze the classification boundary for better feature learning.

3.4. Uncertainty-aware Attention Modeling

As each of the channels looks into different aspects of the spoofing characteristics, we seek to combine those independent channels adaptively to boost the final spoofing detection performance. We introduce an uncertainty-driven attention module that is orthogonal to each of the main and proxy tasks, which thus effectively overcomes the over-confident issue of the traditional attention modules.

Given an input \mathbf{x}_v , the joint likelihood $p(y|\mathbf{x}_v) = p(\mathbf{z}|\mathbf{x}_v)p(\mathbf{d}|\mathbf{x}_v)p(\mathbf{r}|\mathbf{x}_v)p(\mathbf{f}|\mathbf{x}_v)$ according to the channel independence assumption, where \mathbf{z} is from Equation 3 as the main task feature, $\mathbf{d} = \Psi_d \circ \Gamma \circ \Phi(\mathbf{x}_v)$ is from Equation 4 as the reflection feature, $\mathbf{r} = \Psi_r \circ \Gamma \circ \Phi(\mathbf{x}_v)$ is from Equation 5 as the depth feature, $\mathbf{f} = \Psi_m \circ \Gamma \circ \Phi(\mathbf{x})$ is from Equation 6 as the material feature. Maximizing the joint likelihood leads to maximizing the summation of each likelihood:

$$\arg \min - \log(p(y|\mathbf{x}_v)) = - \sum_{\mathbf{u}=\mathbf{z},\mathbf{d},\mathbf{r},\mathbf{f}} \log(p(\mathbf{u}|\mathbf{x}_v)) \quad (9)$$

For each channel, we assume a Gaussian distribution $p(\mathbf{u}|\mathbf{x}_v) \sim \mathcal{N}(\mu_{\mathbf{u}}, \sigma_{\mathbf{u}})$ to capture the uncertainty, where $\mu_{\mathbf{u}}$ is the corresponding channel \mathbf{u} classifier’s separation hyper-plane vector or the mean depth map for the depth channel.

5 $\left\{ \begin{array}{l} p(\mathbf{z}|\mathbf{x}_v) \rightarrow \text{main task (liveness)} \\ p(\mathbf{d}|\mathbf{x}_v) \rightarrow \text{reflection task} \\ p(\mathbf{r}|\mathbf{x}_v) \rightarrow \text{depth task} \\ p(\mathbf{f}|\mathbf{x}_v) \rightarrow \text{material featur task} \end{array} \right.$ (maximize)

Under the probabilistic setting, such $\mu_{\mathbf{u}}$ conforms to another Gaussian distribution $\mathcal{N}(\mu_{\mathbf{u}}, \sigma_{\mu_{\mathbf{u}}})$, where $\sigma_{\mu_{\mathbf{u}}}$ is estimated upon sampling from multiple rounds of training. $\sigma_{\mathbf{u}}$ is independently learned via a two FC-layer network structure in parallel to the feature \mathbf{u} . $\mu_{\mathbf{u}}$ is jointly learned with \mathbf{u} during training and is fixed during the uncertainty training. The objective to learn $\sigma_{\mathbf{u}}$ is defined in the following:

$$\mathcal{L}_{\sigma_{\mathbf{u}}} = \sum_{\mathbf{u}=\mathbf{z},\mathbf{d},\mathbf{r},\mathbf{f}} \left(\frac{\|\mathbf{u} - \mu_{\mathbf{u}}\|^2}{2(\sigma_{\mathbf{u}}^2 + \sigma_{\mu_{\mathbf{u}}}^2)} + \frac{D}{2} \log(\sigma_{\mathbf{u}}^2 + \sigma_{\mu_{\mathbf{u}}}^2) \right) \quad (10)$$

where D is the feature dimension. It is independently optimized after the network is converged. During inference, the network outputs \mathbf{u} and $\sigma_{\mathbf{u}}$ simultaneously. Given $\mu_{\mathbf{u}}$ and $\sigma_{\mu_{\mathbf{u}}}$, we then obtain the uncertainty estimate for each channel with Equation 10.

To sum with, we propose a two-stage training procedure. The first stage consists of the training of liveness main task, proxy tasks and pretext task. And the loss is defined as following:

$$\mathcal{L} = \lambda_v \mathcal{L}_v + \lambda_d \mathcal{L}_d + \lambda_r \mathcal{L}_r + \lambda_m \mathcal{L}_m + \lambda_t (\mathcal{L}_{tid} + \mathcal{L}_{ta}) \quad (11)$$

Then in the second stage, the uncertainty attention module is trained with Equation 10.

4. Implementation Details

In our implementation, we leverage a pre-trained face recognition engine and re-utilize the encoder as our shared feature extractor Φ . Then, we keep the face recognition as our pretext task and equip the main and proxy tasks to form a multi-source multi-channel training. As illustrated in the methodology section, a two-stage training is conducted. For the first stage joint training of pretext, main and proxy tasks, we apply random cropping and horizontal flipping as data augmentation.

We adopt Adam solver and the initial learning rate is set 0.0001. The momentum and weight decay are fixed as 0.9 and 0, respectively. Hyper-parameters in Equation 11 is empirically searched via some hold-out validation as $\lambda_r = 1, \lambda_d = 1, \lambda_v = 1, \lambda_m = 0.1, \lambda_t = 0.1$ for triplet need to add here together with m_1 and m_2 as in Equation 7 and 8, respectively. For the second stage, when training the uncertainty-aware attention module, we re-use well-trained modules from the first stage, and only fine-tune the two-layer fully connected layers for each of the main and proxy tasks to estimate the variance.

5. Experiments

5.1. Datasets

CASIA [60]: A video based 2D spoofing attack database, consists of 600 videos from 50 people, 240 videos from 20

people for training and 360 videos from 30 people for testing. Each people contains 12 videos with video re-display and photo print attacks, of which 8 are normal resolution videos and 4 are high resolution videos. The photo attacks are further categorized into cut photo by cutting holes around eyes, noise, mouth, and warped photos by warping photos with different curvature.

Replay-Attack [13]: A 2D face spoofing attack database consists of 1,300 video clips of photo and video attack attempts from 50 clients, under different lighting conditions. To produce the attacks, high-resolution photos and videos from each client were taken under the same conditions as in their authentication step.

MSU-MFSD [53]: it consists of 280 video clips of photo and video attack attempts from 35 clients. Mobile phones are used to capture both genuine faces and spoofing attacks. Printed photos are generated from high quality color printers for another attack type. It also provides replay video attacks with high resolution of 2048×1536 from iPad air screen.

Oulu-NPU [11]: A large-scale 2D spoofing attack dataset, consists of 4950 genuine and attack videos from 55 people. They are recorded using the front cameras of different mobile devices with variant illuminations and backgrounds. The attack types are print and video replay. There are four protocols designed to consider generalization on cross attack types and capturing sensor types.

SiW [34]: Spoofing in the Wild dataset provides 4,478 genuine and spoofing videos from 165 subjects. For each subject, 8 genuine and up to 20 spoofing videos are captured. It systematically considers variations from subjects, camera sensors, spoofing attack types, lighting conditions, image resolution, and different sessions for capturing. There are three evaluation protocols emphasizing the generalization on face PAD, cross attack types, and unknown attack types.

5.2. Evaluation Metrics

The evaluation is focused on testing the generalization of cross attack types within one dataset, termed intra-dataset evaluation, and cross dataset spoofing detection, termed inter-dataset evaluation following [2]. To be consistent with most of the previous spoofing detection works, we apply the evaluation metrics as: Attack Presentation Classification Error Rate (APCER[19]), Bona Fide Presentation Classification Error Rate (BPCER[19]), ACER = $0.5(\text{APCER} + \text{BPCER})$, Area Under Curve (AUC) ratio and Half-Total Error Rate (HTER). Further following [34] in SiW protocol settings, we use Equal Error Rate (EER) as validation metric for all models to search the threshold to report performance.

5.3. Intra-Dataset Evaluation

We evaluate on a recent large scale spoofing dataset SiW with carefully designed cross attack type testing protocols.

Protocol	1			2			3		
	APCER	BPCER	ACER	APCER	BPCER	ACER	APCER	BPCER	ACER
Methods									
SVM _{RBF} +LBP[11]	4.17	4.17	4.17	5.29±4.39	5.29±4.39	5.29±4.39	16.84±1.89	16.84±1.89	16.84±1.89
SVM _{RBF} +BSIF [2]	7.95	7.95	7.95	7.34±3.30	7.34±3.30	7.34±3.30	25.56±5.63	25.56±5.63	25.56±5.63
FAS-BAS[34]	3.58	3.58	3.58	0.57±0.69	0.57±0.69	0.57±0.69	8.31±3.81	8.31±3.81	8.31±3.81
FAS-TD-SF[52]	1.27	0.83	1.05	0.33±0.27	0.29±0.39	0.31±0.28	7.70±3.88	7.76±4.09	7.73±3.99
Ours (L)	0.66	0.66	0.66	0.35±0.32	0.35±0.32	0.35±0.32	7.97±5.03	7.97±5.03	7.97±5.03
Ours (L+D)	0.47	0.47	0.47	0.25±0.21	0.25±0.21	0.25±0.21	7.75±4.97	7.75±4.97	7.75±4.97
Ours (L+M)	0.57	0.57	0.57	0.27±0.24	0.27±0.24	0.27±0.24	7.80±4.95	7.80±4.95	7.80±4.95
Ours (L+D+R)	0.42	0.42	0.42	0.27±0.22	0.27±0.22	0.27±0.22	7.73±4.96	7.73±4.96	7.73±4.96
Ours (L+D+R+M)	0.44	0.44	0.44	0.24±0.22	0.24±0.22	0.24±0.22	7.52±4.91	7.52±4.91	7.52±4.91
Ours (L+M+A) w/o Triplet	0.56	0.56	0.56	0.27±0.20	0.27±0.20	0.27±0.20	7.773±9.91	7.773±9.91	7.773±9.91
Ours (L+D+A)	0.46	0.46	0.46	0.25±0.22	0.25±0.22	0.25±0.22	7.50±4.79	7.50±4.79	7.50±4.79
Ours (L+D+R+A)	0.40	0.40	0.40	0.23±0.21	0.23±0.21	0.23±0.21	7.43±4.82	7.43±4.82	7.43±4.82
Ours (L+D+R+M+A)	0.36	0.36	0.36	0.20±0.16	0.20±0.16	0.20±0.16	7.32±4.80	7.32±4.80	7.32±4.80

Table 2: Intra-dataset evaluation on SiW dataset. L: main spoofing/liveness detection. L+D: spoofing and depth channels. L+M: spoofing and material channels, L+D+M: spoofing, depth and material channels. L+D+R+M: spoofing, depth, reflection and material channels. L+D+A: spoofing and depth with attention. L+D+M+A: spoofing, depth and material with attention. L+D+R+M+A: our overall model with attention.

Dataset	CASIA				Replay Attack				MSU				All	
	V	C-P	W-P	All	V	D-P	P-P	All	P-P	H-V	M-V	All	Mean	Std
OC-SVM _{RBF} +IMQ[2]	63.26	59.43	66.81	63.34	84.48	67.57	70.30	74.49	53.94	84.75	76.56	72.61	70.14	4.87
OC-SVM _{RBF} +BSIF[2]	67.59	51.01	96.33	72.76	46.54	63.24	38.88	50.62	62.06	80.56	64.06	69.25	64.21	9.71
SVM _{RBF} +LBP[11]	77.41	87.14	69.48	77.61	69.64	73.31	71.85	71.58	55.39	96.02	94.88	83.36	77.51	4.80
NN+LBP[54]	71.80	70.26	67.55	69.78	36.93	75.43	69.45	59.75	26.10	96.84	85.31	71.48	69.75	8.30
GMM+LBP[54]	65.41	85.00	50.15	66.06	60.78	61.46	55.32	59.57	59.35	91.18	86.43	79.92	68.51	8.48
OC-SVM _{RBF} [54]	64.94	85.75	55.15	67.95	84.83	72.62	57.34	73.01	60.90	68.41	75.51	68.60	69.85	2.24
AE+LBP[54]	77.72	80.30	52.92	69.56	79.67	54.92	52.71	63.39	55.67	87.94	92.18	79.67	70.87	6.71
Auxiliary [34]	-	-	-	73.15	-	-	-	71.69	-	-	-	85.88	76.90	6.37
*MADDG [46]	-	-	-	84.51	-	-	-	84.99	-	-	-	88.06	85.85	1.57
Ours (L)	74.88	77.44	81.17	79.31	82.09	72.96	91.42	84.44	66.25	96.60	95.34	85.81	83.18	2.79
Ours (L+D+A)	80.02	81.79	87.80	87.61	85.54	84.37	95.65	84.82	67.82	97.52	96.16	88.59	87.00	1.59
Ours (L+D+R+A)	80.43	81.34	89.24	87.80	86.15	84.56	95.63	85.23	68.14	97.54	97.02	88.24	87.09	1.32
Ours (L+D+R+M+A)	80.69	82.13	90.06	87.92	86.69	84.92	96.33	85.27	68.23	97.70	97.50	89.23	87.47	1.64

Table 3: Inter-dataset evaluation on CASIA, Replay Attack and MSU, AUC(%) is reported. We follow the “Leave one dataset & attack-type out” protocol in [2], where the attack types in testing set is unseen in the training set. We abbreviate V, C-P, W-P, D-P, P-P, H-V and M-V for Video, Cut Photo, Warped Photo, Digital Photo, Printed Photo, HR Video and Mobile video, respectively. *: retrained by their released codes.

We refer another intra-dataset evaluation on Replay-Attack to supplementary due to space limit.

SiW Evaluation: There are 3 protocols in SiW. Protocol 1 focus on evaluating the performance of variations in face pose and expression. Protocol 2 focuses on the unseen medium of replay attack. It chooses 3 out of 4 display attacks, as training and leaving the remaining one as testing, which is iteratively conducted 4 times and averaged. Protocol 3 evaluates cross presentation attack detection, i.e., from print attack to replay attack and vice versa. Averaging over the two is reported.

In Table 2, our method consistently outperforms the other methods with significant margin, i.e., on Protocol 1, we achieve **0.36** ACER while FAS-TD-SF [52] is 1.05. On Protocol 2, ours is **0.20** while the best compared method is 0.31 from FAS-TD-SF. On Protocol 3, ours is **7.32** while the best compared method is 7.73. Similar to Replay-Attack, we apply a gradually increasing module way to highlight effectiveness of the proposed modules. The ablation over our proposed modules suggests: (1) Depth, reflection and mate-

rial are beneficial proxy tasks. (2) putting more proxy tasks together boosts the performance. (3) Our uncertainty-aware attention on top of the baselines can further achieve performance gain with significant margin. We also show an ablation contrasting w/ or w/o using attributed-assisted triplet mining, it shows that by adding triplet constraint, there is continuous margin gain over the other baselines.

5.4. Inter-Dataset Evaluation

The inter-dataset setting mimics the real setting for unseen attack across types and datasets. We consider two protocols. One follows the traditional rule [2], a “Leave one dataset & attack-type out” protocol, taking CASIA, Replay-Attack and MSU-MFSD as our datasets. Each of the three datasets contains three attack types. When evaluating one attack type of one dataset, we pick the other two datasets for training and excluding the testing attack type from training. The other less-strict setting is the “Leave one dataset out” protocol used in MADDG[46], the difference to the former is that in this protocol, training and testing sets would have

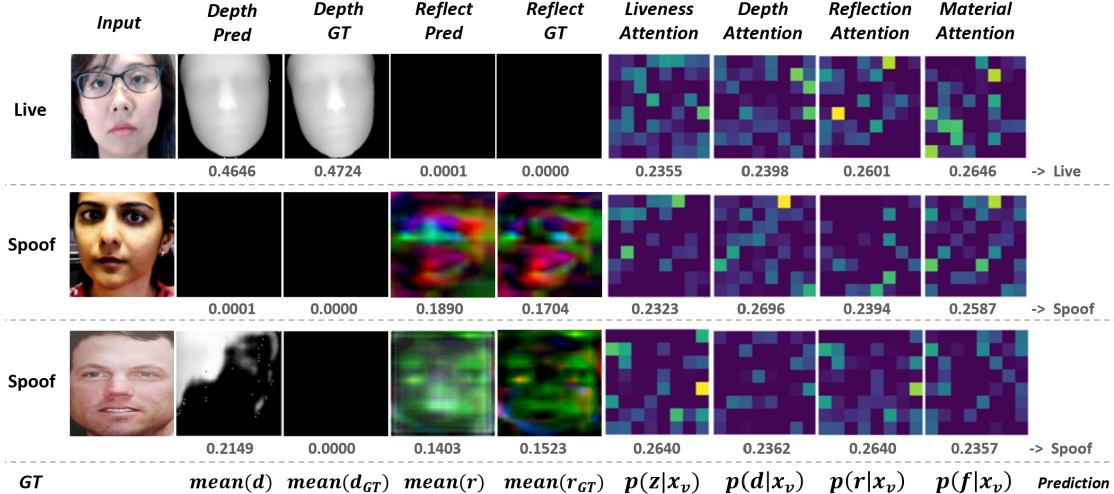


Figure 4: Visual examples of our Uncertainty-aware Attention Modeling. We show example of Spoof and Live Image on SiW dataset. The first two rows shows examples is from SiW dataset while the bottom row is from CelebA-Spoof dataset [59].

Method	DD	OCI-M	OMI-C	OCM-I	ICM-O
MS LBP[37]		78.50	44.98	51.64	49.31
Binary CNN[55]		82.87	71.94	65.88	77.54
IDA [53]	×	66.67	39.05	78.25	44.59
Color Texture [10]		78.47	76.89	62.78	32.71
LBPTOP [14]		70.80	61.05	49.54	44.09
Auxiliary [34]		85.88	73.15	71.69	77.61
Ours (L)		85.81	79.31	84.44	76.26
Ours (L+D+A)	×	88.59	87.61	86.74	80.30
Ours (L+D+R+A)		90.73	89.03	88.42	83.58
Ours (L+D+R+M+A)		91.32	89.28	91.83	85.48
MADDG[47]		88.06	84.51	84.99	80.02
RFGML [48]	✓	93.98	88.16	90.48	91.16
SSDG-M [20]		90.47	85.45	94.61	81.83

Table 4: Inter-dataset evaluation on CASIA, Replay Attack, MSU and Oulu-NPU dataset. AUC (%) is reported. We follow the “Leave one dataset out” protocol in [46], where training and testing sets share attack types. DD denotes disentangling source domains.

overlapping attack types.

Leave one dataset & attack-type out In Table 3, we evaluate each of the three attack types from three datasets. Both traditional feature learning based methods and most recent deep learning based methods [34, 46] are compared. Overall we achieve consistently stronger results than the other methods. In CASIA, video attack is significantly better than other methods while cut photo and warped photo are among the top. In Replay-Attack, we achieve clear better performance. In MSU-MFSD, we observe 1% to 6% performance improvement over the compared methods.

Leave one dataset out In Table 4. Since our physically-guided proxy task does not require any domain priors, we compare methods w and w/o source domains disentanglement (DD). Our method surpass all methods without domain disentanglement including [34], while still achieve

comparable performance compare to methods [47][48][20] that utilize extra source domains information.

5.5. Analysis of Uncertainty-aware Attention

We visualize our Uncertainty-aware Attention Module in Figure 4. Specifically, we visualize the last feature map of the last FC layer in the Attention Module across Liveness, Depth, Reflection and Spoof channel alongside with the corresponding Input Image, GT/Predicted Depth map and GT/Predicted Reflection map. In Figure 4, we can see that, in most cases, those proxy tasks agree with each other like the examples from top two rows, while in some cases, those proxy tasks does not agree with each other. for example, the figure in the bottom row, depth channel made the wrong prediction and disagree with other channels, However, our uncertainty-aware attention module are able to correct his by give depth a lower confidence and voted for prediction from other channels, thus correcting the final decision.

6. Conclusion

In this work, we propose depth, reflection and material guided proxy tasks for unseen spoofing attacks. We propose a multi-source multi-channel training scheme for model optimization. Due to the consistency of depth, reflection and skin material across different spoofing scenario on genuine faces, by harnessing those physical proxy tasks, we expect the proposed method to deal with unseen spoofing attacks. Finally, an uncertainty-aware attention module is introduced to aggregate the multiple channels for boosted performance. Experiments across intra- and inter-dataset protocols show our method achieves consistently better performance and is effective for unseen spoofing detection.

References

- [1] A. Agarwal, R. Singh, , and M. Vatsa. Face anti-spoofing using haralick features. In *BATS*, 2016. 3
- [2] S.R. Arashloo, J. Kittler, and W. Christmas. Anomaly detection approach to face spoofing detection: a new formulation and evaluation protocol. *IEEE Access*, 2017. 3, 6, 7
- [3] Y. Atoum, Y. Liu, A. Jourabloo, and X. Liu. Face anti-spoofing using patch and depth-based cnns. In *IJCB*, 2017. 3
- [4] W. Bao, H. Li, N. Li, and W. Jiang. A liveness detection method for face recognition based on optical flow field. In *IASP*, 2009. 3
- [5] Sean Bell, Paul Upchurch, Noah Snaveley, and Kavita Bala. Material recognition in the wild with the materials in context database. In *CVPR*, 2015. 1, 2, 5
- [6] G Betta, D Capriglione, C Liguori, and A Paolillo. Uncertainty evaluation in face recognition algorithms. In *2011 IEEE International Instrumentation and Measurement Technology Conference*, pages 1–6. IEEE, 2011. 3
- [7] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *SIGGRAPH*, 1999. 1
- [8] S. Bobbia, Y. Benezeth, and J. Dubois. Remote photoplethysmography based on implicit living skin tissue segmentation. In *ICPR*, 2016. 2
- [9] Z. Boulkenafet, J. Komulainen, and A. Hadid. Face anti-spoofing based on color texture analysis. In *ICIP*, 2015. 3
- [10] 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. 8
- [11] Zinelabinde Boulkenafet, Jukka Komulainen, Lei Li, Xiaoyi Feng, and Abdenour Hadid. Oulu-npu: A mobile face presentation attack database with real-world variations. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*, pages 612–618. IEEE, 2017. 6, 7
- [12] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 4
- [13] Ivana Chingovska, André Anjos, and Sébastien 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. IEEE, 2012. 6
- [14] 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. 8
- [15] G. de Haan and V. Jeanne. Robust pulse rate from chrominance-based rppg. *IEEE Trans. Biomedical Engineering*, 2013. 2
- [16] L. Feng, L. Po, Y. Li, X. Xu, F. Yuan, T.C. Cheung, and K. Cheung. Integration of image quality and motion cues for face anti-spoofing: A neural network approach. *Journal of Visual Communication and Image Representation*, 2016. 3
- [17] Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *ECCV*, 2018. 2, 4
- [18] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *CVPR*, 2014. 3
- [19] ISO Central Secretary. Information technology — biometric presentation attack detection — part 1: Framework. Standard ISO/IEC 30107-1:2016, International Organization for Standardization, Geneva, CH, 2016. 6
- [20] Yunpei Jia, Jie Zhang, Shiguang Shan, and Xilin Chen. Single-side domain generalization for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8484–8493, 2020. 8
- [21] Ajay J Joshi, Fatih Porikli, and Nikolaos Papanikolopoulos. Multi-class active learning for image classification. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2372–2379. IEEE, 2009. 3
- [22] A. Jourabloo, Y. Liu, and X. Liu. Face de-spoofing: Anti-spoofing via noise modeling. In *ECCV*, 2018. 3
- [23] J. Yang, Z. Lei, S. Liao, and S. Z. Li. Face liveness detection with component dependent descriptor. In *ICB*, 2013. 3
- [24] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017. 3
- [25] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018. 3
- [26] Taewook Kim, YongHyun Kim, Inhan Kim, and Daijin Kim. Basn: Enriching feature representation using bipartite auxiliary supervisions for face anti-spoofing. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2
- [27] K. Kollreider, H. Fronthaler, M. I. Faraj, and J. Bigun. Real-time face detection and motion analysis with application in liveness assessment. *TIFS*, 2007. 2
- [28] J. Komulainen, A. Hadid, and M. Pietikainen. Context based face anti-spoofing. In *BATS*, 2013. 3
- [29] C.Y. Lee, S. Xie, P.W. Gallagher, Z. Zhang, and Z. Tu. Deeply-supervised nets. In *AISTATS*, 2015. 3
- [30] David D Lewis and Jason Catlett. Heterogeneous uncertainty sampling for supervised learning. In *Machine learning proceedings 1994*, pages 148–156. Elsevier, 1994. 3
- [31] J. Li, Y. Wang, T. Tan, and A. K. Jain. Live face detection based on the analysis of fourier spectra. In *SPIE*, 2004. 3
- [32] L. Li, X. Feng, Z. Boulkenafet, Z. Xia, M. Li, and A. Hadid. An original face anti-spoofing approach using partial convolutional neural network. In *IPTA*, 2016. 3
- [33] S. Liu, P.C. Yuen, and S. Zhang G. Zhao. 3d mask face anti-spoofing with remote photoplethysmography. In *ECCV*, 2016. 2
- [34] Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. Learning deep models for face anti-spoofing: Binary or auxiliary supervision. In *CVPR*, 2018. 2, 3, 6, 7, 8

- [35] Y. Liu, J. Stehouwer, A. Jourabloo, and X. Liu. Deep tree learning for zero-shot face antispoofing. In *CVPR*, 2019. 3
- [36] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, 2015. 5
- [37] Jukka Määttä, Abdenour Hadid, and Matti Pietikäinen. Face spoofing detection from single images using micro-texture analysis. In *2011 international joint conference on Biometrics (IJCB)*, pages 1–7. IEEE, 2011. 8
- [38] Jishnu Mukhoti and Yarin Gal. Evaluating bayesian deep learning methods for semantic segmentation. *arXiv preprint arXiv:1811.12709*, 2018. 3
- [39] Jukka Määttä, Abdenour Hadid, and Matti Pietikäinen. Face spoofing detection from single images using micro-texture analysis. In *IJCB*, 2011. 3
- [40] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-glass networks for human pose estimation. In *ECCV*, 2016. 4
- [41] E. M. Nowara, A. Sabharwal, and A. Veeraraghavan. Ppgsecure: Biometric presentation attack detection using photoplethysmograms. In *FG*, 2017. 2
- [42] G. Pan, L. Sun, Z. Wu, and S. Lao. Eyeblink based anti-spoofing in face recognition from a generic webcam. In *ICCV*, 2007. 2
- [43] K. Patel, H. Han, and A.K. Jain. Secure face unlock: Spoof detection on smartphones. *TIFS*, 2016. 3
- [44] K. Patel, H. Han, and A. K. Jain. Cross-database face anti-spoofing with robust feature representation. In *CCBR*, 2016. 3
- [45] T. Pereira, A. Anjos, J.M. DeMartino, and S. Marcel. Lbp-top based counter measure against face spoofing attacks. In *ACCV*, 2012. 3
- [46] Rui Shao, Xiangyuan Lan, Jiawei Li, and Pong C. Yuen. Multi-adversarial discriminative deep domain generalization for face presentation attack detection. In *CVPR*, 2019. 2, 3, 7, 8
- [47] Rui Shao, Xiangyuan Lan, Jiawei Li, and Pong C Yuen. Multi-adversarial discriminative deep domain generalization for face presentation attack detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10023–10031, 2019. 8
- [48] Rui Shao, Xiangyuan Lan, and Pong C Yuen. Regularized fine-grained meta face anti-spoofing. In *AAAI*, pages 11974–11981, 2020. 8
- [49] T.A. Siddiqui, S. Bharadwaj, T.I. Dhamecha, A. Agarwal, M. Vatsa, R. Singh, and N. Ratha. Face anti-spoofing with multifeature videolet aggregation. In *ICPR*, 2016. 3
- [50] Qing Sun, Ankit Laddha, and Dhruv Batra. Active learning for structured probabilistic models with histogram approximation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3612–3621, 2015. 3
- [51] X. Tan, Y. Li, J. Liu, and L. Jiang. Face liveness detection from a single image with sparse low rank bilinear discriminative model. In *ECCV*, 2010. 3
- [52] Zezheng Wang, Chenxu Zhao, Yunxiao Qin, Qiusheng Zhou, and Zhen Lei. Exploiting temporal and depth information for multi-frame face anti-spoofing. *arXiv preprint arXiv:1811.05118*, 2018. 3, 7
- [53] 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. 6, 8
- [54] Fei Xiong and Wael AbdAlmageed. Unknown presentation attack detection with face rgb images. In *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, pages 1–9. IEEE, 2018. 3, 7
- [55] Jianwei Yang, Zhen Lei, and Stan Z Li. Learn convolutional neural network for face anti-spoofing. *arXiv preprint arXiv:1408.5601*, 2014. 8
- [56] Z. Yang, Z. Lei, and S.Z. Li. Learn convolutional neural network for face anti-spoofing. In *arXiv:1408.5601*, 2014. 3
- [57] 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 *CVPR*, 2019. 3
- [58] Xuaner Zhang, Ren Ng, and Qifeng Chen. Single image reflection separation with perceptual losses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4786–4794, 2018. 1, 2, 4
- [59] Yuanhan Zhang, ZhenFei Yin, Yidong Li, Guojun Yin, Junjie Yan, Jing Shao, and Ziwei Liu. Celeba-spoof: Large-scale face anti-spoofing dataset with rich annotations. In *European Conference on Computer Vision*, pages 70–85. Springer, 2020. 5, 8
- [60] Zhiwei Zhang, Junjie Yan, Sifei Liu, Zhen Lei, Dong Yi, and Stan Z Li. A face antispoofing database with diverse attacks. In *2012 5th IAPR international conference on Biometrics (ICB)*, pages 26–31. IEEE, 2012. 6
- [61] Z.Xu, S.Li, and W.Deng. Learning temporal features using lstm-cnn architecture for face anti-spoofing. In *ACPR*, 2016. 3