

Learning Generalized Spoof Cues for Face Anti-spoofing

Haocheng Feng¹, Zhibin Hong¹, Haixiao Yue¹, Yang Chen^{2*}, Keyao Wang¹,
Junyu Han¹, Jingtuo Liu¹, and Errui Ding¹

Department of Computer Vision Technology(VIS), Baidu Inc.¹
Beihang University²

{fenghaocheng, hongzhibin, yuehaixiao, wangkeyao,
hanjunyu, liujingtuo, dingerrui}@baidu.com
chenyangwiz@buaa.edu.cn

Abstract. Many existing face anti-spoofing (FAS) methods focus on modeling the decision boundaries for some predefined spoof types. However, the diversity of the spoof samples including the unknown ones hinders the effective decision boundary modeling and leads to weak generalization capability. In this paper, we reformulate FAS in an anomaly detection perspective and propose a residual-learning framework to learn the discriminative live-spoof differences which are defined as the spoof cues. The proposed framework consists of a spoof cue generator and an auxiliary classifier. The generator minimizes the spoof cues of live samples while imposes no explicit constraint on those of spoof samples to generalize well to unseen attacks. In this way, anomaly detection is implicitly used to guide spoof cue generation, leading to discriminative feature learning. The auxiliary classifier serves as a spoof cue amplifier and makes the spoof cues more discriminative. We conduct extensive experiments and the experimental results show the proposed method consistently outperforms the state-of-the-art methods. The code will be publicly available at <https://github.com/vis-var/lgsc-for-fas>.

Keywords: Face anti-spoofing, Anomaly detection, Residual learning

1 Introduction

Face authentication applications have been widely used in our daily lives. However, hackers may fraud the face recognition systems by the means of presentation attack (PA) such as printed photos (*i.e.* print attack), digital images / videos (*i.e.* replay attack) and 3D facial masks (*i.e.* 3D mask attack). To this end, face anti-spoofing (FAS) is developed to guarantee the security of the face recognition system.

Traditional methods employ hand-crafted features such as LBP [24], HOG [19] and SIFT [27] to extract the texture information and use shallow classifiers

* This work was done when Yang Chen was an intern at Baidu Inc.

to model the decision boundary. The hand-crafted features used by such methods are not discriminative enough and the classifiers performance is limited. To take advantage of better feature representations, the early deep learning based methods [34,35] utilize CNN to learn more discriminative features and formulate the spoofing detection as a binary classification problem. However, these methods tend to overfit on the predefined datasets and cannot generalize well. Most recent deep learning based methods [22,30,32] focus on improving the model’s generalization capacity and achieve promising results.

The limited generalization capacity of FAS lies in the diversity of the spoof samples including the unknown ones. Though one can assume that the live samples share the same nature and can be categorized into one class, the spoof samples are generally diverse due to the wide variety of attack mediums. To relieve the spoof diversity’s influence on the decision boundary modeling, we reformulate FAS from an anomaly detection perspective and assume the live samples belong to a closed-set while the spoof ones are outliers of the closed-set and belong to an open-set. Based on the hypothesis, we define the spoof cues as the discriminative features that can be used to differentiate the closed-set and the open-set.

In this paper, the spoof cues are formulated as a feature map that has the same size as the input image. The spoof cue map is prone to be a non-zero map for a spoof image while an all-zero map for a live image. Unlike the compact embeddings which are traditionally used for classification (live vs. spoof), spoof cues can effectively encode the spatial information. To learn the spoof cue map, we propose the residual-learning framework that consists of a spoof cue generator and an auxiliary classifier. In the spoof cue generator, we set explicit regression loss for live samples to minimize the magnitudes of their spoof cues while setting no explicit constraint on the spoof samples, leading their corresponding entries to any real numbers. In this way, the spoof cues of live and spoof samples are naturally separable. In addition, we adopt the multi-scale feature-level metric learning on the generator to promote live-live intra-class compactness and live-spoof inter-class separability. Then, the spoof cue map and the input image, which are skip-connected in a residual learning manner, are fed to an auxiliary classifier, further improving the discriminative feature learning. The spoof cue maps and their distributions are shown in Fig. 1. The spoof cue maps are directly utilized to calculate the score of an input sample being spoof in the test phase. The experimental results show that the spoof cue maps can be utilized to separate the live closed-set and spoof open-set effectively.

Our contributions can be summarized as:

1. We reformulate FAS from an anomaly detection perspective and propose a spoof cue generator, which models the spoof cues of live samples on closed-set and spoof samples on open-set by minimizing the magnitudes of spoof cues for live samples while imposing no explicit constraint on the spoof samples.
2. We deploy a residual learning framework to make the spoof cues more discriminative and further boost the spoof cues’ generalization capability.

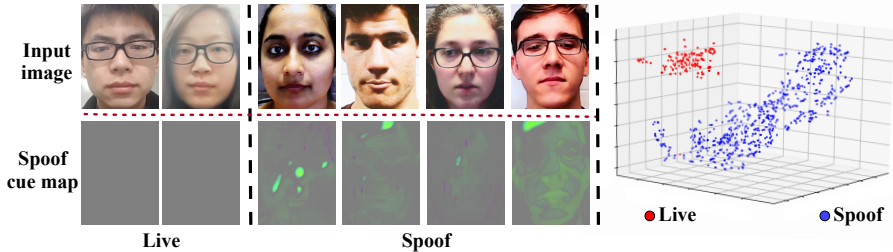


Fig. 1. The illustration of the spoof cue maps and the distribution of the learned feature embeddings. **Left:** the input images and the generated spoof cue maps. **Right:** the distribution of the learned feature embeddings visualized using the t-SNE. More details will be given in Section 4.4.

3. Without additional depth or temporal information, the proposed framework outperforms the state-of-the-art on the popular RGB anti-spoofing datasets.

2 Related Work

2.1 Face anti-spoofing methods

The traditional face anti-spoofing methods use the hand-crafted features, such as LBP [24], HOG [19] and SIFT [27] to represent the image texture and use the traditional classifiers like SVM for binary classification. Besides, some works also adopt the feature extractors in the different color spaces like HSV and YCbCr [6] to tackle face anti-spoofing. However, the hand-crafted features are specially designed for several defined types of spoof and are not robust enough to variable conditions such as the illumination and the presentation attack types. Motion cues like eye-blinking [26] and head movement [18] are also used when the temporal information is given. However, the performance of these methods drop drastically on the video replay attack.

With the development of CNN, many researchers start to use CNN as a feature extractor to exploit the discriminative representations in RGB images and formulate face anti-spoofing as a binary classification problem [34]. Moreover, Feng *et al.* [11] utilize LSTM in temporal information modeling for better classification performance. Despite their good performances in the intra-dataset test, these methods are easy to overfit and have poor generalizability.

Most recent methods can be categorized into three classes: the additional information-based methods, the domain shift methods, and the few-shot methods. These methods are designed to further boost the models' generalizability. Liu *et al.* [22] use pseudo-depth to supervise the model training. Shao *et al.* [32] cast FAS as a domain generalization problem and develop a Regularized Fine-grained Meta-learning framework. Qin *et al.* [30] propose a novel Adaptive Inner-update Meta Face Anti-Spoofing (AIM-FAS) method.

2.2 Anomaly detection

Anomaly detection is the task to identify the unusual samples from a set of normal data. It has been widely used in a variety of applications such as fraud detection [1], cyber-intrusion detection [20], health care [33] and video surveillance [4]. In the face anti-spoofing area, Liu *et al.* [3] formulate the face PAD problem as anomaly detection and put forward a one-class classifier. Nikisins *et al.* [25] use a Gaussian Mixture model-based anomaly detector. Arashloo *et al.* [2] use client-specific information in a one-class anomaly detection formulation to improve the model’s performance significantly. Perez *et al.* [28] reformat the FAS problem from an anomaly detection perspective and put forward a deep metric learning model with triplet focal loss.

3 Our Approach

3.1 Preliminaries

Anomaly detection. Anomaly detection (AD) is the task of identifying the unusual samples in a set of normal data. The typical AD methods attempt to learn a compact description of the normal data in an unsupervised manner. Take the deep SVVD [31] as an example, for input space ($\mathcal{X} \subseteq \mathbb{R}^d$) and output space ($\mathcal{Z} \subseteq \mathbb{R}^p$), let $\phi(\cdot; \mathcal{W}) : \mathcal{X} \rightarrow \mathcal{Z}$ be a neural network with L hidden layers and corresponding set of weights $\mathcal{W} = \{\mathcal{W}^1, \dots, \mathcal{W}^L\}$. The objective is to train the neural network ϕ to learn a transformation that minimizes the volume of a data-enclosing hypersphere in output space \mathcal{Z} centered on a predetermined point c . Given N (unlabeled) training samples $(x_1, \dots, x_N \subseteq \mathcal{X})$, the objective is:

$$\min_{\mathcal{W}} \frac{1}{N} \sum_{i=1}^N \|\phi(x_i; \mathcal{W}) - c\|^2 + \frac{\lambda}{2} \sum_{\ell=1}^L \|\mathcal{W}^\ell\|_F^2, \lambda > 0, \quad (1)$$

Once the network is trained, the anomaly score for a test point x is given by the distance from $\phi(x; \mathcal{W})$ to the center of the hypersphere:

$$s(x) = \|\phi(x; \mathcal{W}) - c\|. \quad (2)$$

In FAS, though the live samples are assumed to share the same nature, the spoof samples can be very diverse due to the wide variety of attack mediums. Such diversity makes it hard for the spoof samples to form a compact region in the feature representation space and further hinders the effective decision boundary modeling between live and spoof samples. Besides, the decision boundary learned by the known spoof samples may have poor performance on the unseen spoof samples. Motivated by the AD method, we assume the live samples belong to a closed-set while the spoof samples are outliers from this closed-set and belong to an open-set. Therefore, the distribution of the live samples is assumed to lie in a compact sphere in the learned feature representation space while the spoof samples far away from the center of the live sphere. Specifically, for the face anti-spoofing model’s input space ($\mathcal{X} \subseteq \mathbb{R}^d$) and output space ($\mathcal{Z} \subseteq \mathbb{R}^p$). Let $\phi(\cdot; \mathcal{W}) : \mathcal{X} \rightarrow \mathcal{Z}$ be a neural network with L hidden layers and corresponding set of weights $\mathcal{W} = \{\mathcal{W}^1, \dots, \mathcal{W}^L\}$. Given N_l live samples $(x_1, \dots, x_{N_l} \subseteq \mathcal{X})$,

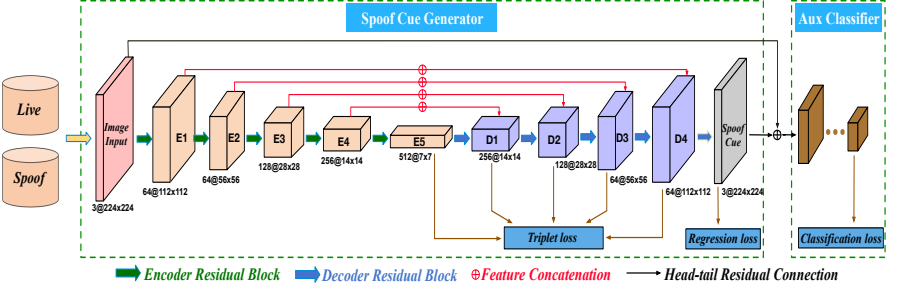


Fig. 2. The proposed network architecture.

N_s spoof samples ($y_1, \dots, y_{N_s} \subseteq \mathcal{X}$), let c be the center of the live samples in the output space \mathcal{Z} , the objective is:

$$\min_{\mathcal{W}} \frac{1}{N_l} \sum_{i=1}^{N_l} \|\phi(x_i; \mathcal{W}) - c\|^2, \quad (3)$$

$$\max_{\mathcal{W}} \frac{1}{N_s} \sum_{i=1}^{N_s} \|\phi(y_i; \mathcal{W}) - c\|^2. \quad (4)$$

In the test phase, the score for a test sample t being spoof can be given by the distance from $\phi(t; \mathcal{W})$ to the center of the live hypersphere:

$$s(x) = \|\phi(t; \mathcal{W}) - c\|. \quad (5)$$

In our proposed method, we utilize the explicit regression supervision on the live samples to achieve the optimization goal of Eq. 3 and the implicit metric learning supervision on the live and spoof samples to achieve the optimization goal of Eq. 4. In the test phase, we directly use the distance between the test samples and the predefined live closed set's center in the feature space to calculate the score of a sample being spoof.

Residual learning. Deep residual learning [12] has gained great success in visual recognition, the residual learning method recasts the desired underlying mapping $H(x)$ into $H(x) = F(x) + x$ and learns the effective residual feature representation.

In FAS, the discriminative differences between the live and spoof samples may be caused by the intrinsic properties of the mediums such as printed papers or screens. Motivated by the residual learning method, the discriminative differences between live and spoof samples can be regarded as the residual of the spoof samples and the live samples. We define the discriminative differences as spoof cues and assume the spoof cues only exist in the spoof samples. In order to learn the spoof cues, we propose a residual learning framework that consists of a spoof cue generator and an auxiliary classifier. In the spoof cue generator, we impose the regression loss on the live samples while put no explicit supervision

on the spoof samples, which guarantees the spoof cues’ generalization capability. The auxiliary classifier serves as a spoof cue amplifier in the end-to-end training process and helps to learn more discriminative spoof cues.

3.2 Spoof cue learning framework

The proposed residual learning framework is illustrated in Fig. 2. It consists of a spoof cue generator and an auxiliary classifier. In our assumption, the live samples share some common properties that spoof samples do not have. Therefore, all the live samples belong to a closed-set while the spoof samples are outliers from this closed-set but belong to an open-set. The discriminative differences between the live and spoof samples which are defined as spoof cues can be caused by the intrinsic properties of the medium carriers such as printed paper or screens. Our framework aims to learn the discriminative spoof cues and further boost the spoof cues’ generalization capability.

In the spoof cue generator, we adopt a simple yet effective U-Net to generate spoof cue maps. In our assumption, the spoof cues only exist in spoof samples. Thus, for live samples, the spoof cue map should be an all-zero map. For spoof samples, the spoof cue map should be an unknown non-zero map. To relieve the spoof diversitys influence on the decision boundary modeling and make the spoof cue generalized enough, we introduce the pixel-wise regression loss as explicit supervision only for the live samples. To promote the intra-class compactness in live samples and the inter-class separability between live and spoof samples, we introduce the multi-scale triplet loss on the feature embeddings.

In the auxiliary classifier, we feed the spoof cue map and the input image into the auxiliary classifier in a residual learning manner. The auxiliary classifier serves as an amplifier in the end-to-end training process and makes the learned spoof cues more discriminative.

3.3 Spoof Cue Generator

Architecture. In the spoof cue generator, as shown in Fig. 2, we adopt the U-Net architecture which builds skip connections from encoder to decoder in multiple scales to generate spoof cues. We choose the ResNet18 [12] pre-trained on ImageNet [10] as encoder E which contains four Encoder Residual Blocks. Followed by E, the decoder D, composed of five Decoder Residual Blocks, decodes the information back to generate the spoof cues. In each Decoder Residual Block, the feature map from the previous layer is upsampled by a nearest-neighbor interpolation, following which we add a 2×2 convolution. We aggregate the feature maps from the symmetric position in the encoder to the decoder by the concatenation operation and add an Encoder Residual Block afterward. At the end of the decoder, a Decoder Residual Block accompanied by a Tanh activation layer is applied to generate the outputs.

Regression loss. In our assumption, the live samples belong to a closed-set while the spoof samples are outliers from this closed-set and belong to an open-set. Therefore, we do not set any preset constraint on the spoof cues and only

impose constraints on the live ones. Since we assume the spoof cues only exist in spoof samples, the spoof cue maps should be zero maps for live samples while remaining unknown for spoof samples. From our perspective, the zero map can be considered as the center of the live samples in the feature space and the live-only regression loss achieves the optimization goal of Eq. 3.

Given an RGB image I as input, the spoof cue generator outputs a spoof cue map of the same size, the spoof cue map C is a zero map for a live sample. While C remains unknown for a spoof sample. The spoof cue regression loss for a live sample is the pixel-wise $L1$ loss, in the formulation of:

$$L_r = \frac{1}{N_l} \sum_{I_i \in \text{live}} \|C_i\|_1, \quad (6)$$

where N_l is the number of live samples in one batch.

Metric Learning. The metric learning-based loss is designed as an implicit supervision for spoof samples to promote live-live intra-class compactness and live-spoof inter-class separability on the feature level, which achieves the optimization goal of Eq. 4. Specifically, we obtain a set of feature vectors $\{V\}$ by employing the global average pooling (GAP) on feature maps from layer E5 to D4 and apply the metric supervision afterward. We use the triplet loss of which the anchor always belongs to the live class. The metric learning-based loss can be formulated as:

$$L_t = \frac{1}{T} \sum_{i=1}^T \max(d(a_i, p_i) - d(a_i, n_i) + m, 0), \quad (7)$$

$$d(i, j) = \left\| \frac{v_i}{\|v_i\|_2} - \frac{v_j}{\|v_j\|_2} \right\|_2,$$

where $\{a_i, p_i, n_i\}$ denotes the anchor (live), positive (live), negative (spoof) samples within the i th triplet respectively, T denotes the number of triplets, $d(i, j)$ represents the euclidean distance between two L2-normalized feature vectors output by the GAP layer, and m is the pre-defined margin constant.

As for the triplet sampling strategy, we choose the online batch-all triplet mining proposed in [14]. At each training step, we collect all the valid triplets within the current batch of data for metric loss computation, the triplets satisfying $\|d(a, n) - d(a, p)\|_2 < m$.

3.4 Auxiliary Classifier

The success of prior work [13] proves the metric learning based loss in feature embeddings coupled with classification supervision helps to learn more discriminative features, which inspires us to design the auxiliary classifier. In our residual learning framework, the auxiliary classifier serves as a spoof cue amplifier and helps to learn more discriminative spoof cues. After the spoof cue generator, the generated spoof cue maps C are added back to the input image I to form the overlaid images S . With S as input, the auxiliary classification loss can be

formulated as:

$$L_a = \frac{1}{N} \sum_{i=1}^N z_i \log q_i + (1 - z_i) \log(1 - q_i), \quad (8)$$

where N is the number of samples, z_i is the binary label and q_i is the network prediction.

In our visualization experiments, we notice that choosing S instead of C as the input of the auxiliary classifier helps to learn more discriminative spoof cue maps, which proves the residual learning’s superiority. We discuss the experimental details in Section 4.4.

3.5 Training and Testing

Loss functions. The loss functions of the proposed model are three-fold: the pixel-wise $L1$ loss L_r for spoof cue regression on live samples, the triplet loss L_t and the auxiliary binary classification loss L_a on both live and spoof samples. We integrate all these losses and establish the total loss of L during training:

$$L = \alpha_1 L_r + \alpha_2 \sum_{k \in \{E5-D4\}} L_t^k + \alpha_3 L_a, \quad (9)$$

where k indexes the layer where we apply the triplet loss, and $\alpha_1, \alpha_2, \alpha_3$ are the weights to balance the influence of the different loss components.

Test Strategy. At the test stage, we use the generated spoof cue map instead of the classifier’s output to evaluate. In Eq. 5, the distance from $\phi(t; \mathcal{W})$ to the center of the live hypersphere is utilized to calculate the spoof score. In our proposed method the all-zero map can be considered as the center of the live samples in the feature space. We directly use the spoof cue map to evaluate. We obtain the spoof cue map \hat{C} and define the spoof score as the element-wise mean of \hat{C} magnitude. The spoof score evaluates the probability of the test sample being spoof:

$$score = \overline{||\hat{C}||_1}, \quad (10)$$

The input image is likely to be a spoof one when the spoof score is high.

4 Experiments

4.1 Experiment Settings

In our experiments, we crop the human faces out of the images and feed them to the proposed framework. For datasets that offer no face location ground truth, we use the Dlib [16] toolbox as the face detector. We resample the training examples to keep the live-spoof ratio to 1:1. The training batch size is 32 and the warm-up learning rate strategy is applied at the first training epoch. We choose Adam [17] as the optimizer with initial learning rate 1e-3, which decays every 600 training steps by a factor of 0.95. In the spoof cue generator, we use Tanh as the activation function for the output layer of the U-Net, so the output

Table 1. The intra testing results and comparison with the state-of-the-art methods on 3 protocols of SiW dataset.

Protocol	Method	APCER (%)	BPCER (%)	ACER (%)
1	Auxiliary [22]	3.58	3.58	3.58
	STASN [36]	–	–	1.00
	Meta-FAS-DR [38]	0.52	0.50	0.51
	Ours	0.00	0.50	0.25
2	Auxiliary [22]	0.57±0.69	0.57±0.69	0.57±0.69
	Meta-FAS-DR [38]	0.25±0.32	0.33±0.27	0.29±0.28
	STASN [36]	–	–	0.28±0.05
	Ours	0.00±0.00	0.00±0.00	0.00±0.00
3	STASN [36]	–	–	12.10±1.50
	Auxiliary [22]	8.31±3.81	8.31±3.80	8.31±3.81
	Meta-FAS-DR [38]	7.98±4.98	7.35±5.67	7.66±5.32
	Ours	1.61±1.69	0.77±1.09	1.19±1.39

distributes in the same range with the input image, which is linearly normalized to $[-1, 1]$. We set m to 0.5 in the triplet loss and set α_1 to α_3 as 5, 1 and 5, respectively. The score threshold in Eq. 10 is set experimentally, and 0.01 is recommended. Our framework can be trained end-to-end and converges in 20 epochs on all the datasets used in the paper.

As for the metrics, we use Average Classification Error Rate (ACER) which is half of the summation of Attack Presentation Classification Error Rate (APCER) and Bona Fide Presentation Classification Error Rate (BPCER) for intra-dataset test on SiW [22] and OULU-NPU [8] datasets:

$$ACER = \frac{APCER + BPCER}{2}. \quad (11)$$

Note that APCER corresponds to the highest false positive rate in all presentation attack instruments.

For inter-dataset test between CASIA-MFSD [37] and Replay-Attack [9] datasets, we use the Half Total Error Rate (HTER) which is the mean of False Rejection Rate (FRR) and False Acceptance Rate (FAR):

$$HTER = \frac{FRR + FAR}{2}. \quad (12)$$

4.2 Intra Testing

We carry out the intra testing on SiW and OULU-NPU datasets. The SiW dataset collects 4478 HD videos from 165 individuals in total. It contains data in the rich PIE (pose, illumination, and expression) variations. There are four pre-defined protocols on the SiW dataset to evaluate the generalization for face anti-spoofing. The results on the test set are reported in Table 1. The state-of-the-art methods we compare with include Auxiliary [22], STASN [36] and Meta-FAS-DR [38]. Our method outperforms the state-of-the-art results on all protocols by a large margin. To be notified, protocol 3 evaluates the model’s

Table 2. The intra testing results and comparison with the state-of-the-art methods on 4 protocols of OULU-NPU dataset.

Protocol	Method	APCER (%)	BPCER (%)	ACER (%)
1	MILHP [21]	8.3	0.8	4.6
	STASN [36]	1.2	2.5	1.9
	Auxiliary [22]	1.6	1.6	1.6
	FaceDs [15]	1.2	1.7	1.5
	Ours	0.8	0.0	0.4
2	FaceDs [15]	4.2	4.4	4.3
	Auxiliary [22]	2.7	2.7	2.7
	GRANDINT [5]	3.1	1.9	2.5
	STASN [36]	4.2	0.3	2.2
	Ours	0.8	0.6	0.7
3	GRADIANT [5]	2.6±3.9	5.0±5.3	3.8±2.4
	FaceDs [15]	4.0±1.8	3.8±1.2	3.6±1.6
	Auxiliary [22]	2.7±1.3	3.1±1.7	2.9±1.5
	STASN [36]	4.7±3.9	0.9±1.2	2.8±1.6
	Ours	1.5±1.4	1.9±1.9	1.7±1.6
4	GRADIANT [5]	5.0±4.5	15.0±7.1	10.0±5.0
	Auxiliary [22]	9.3±5.6	10.4±6.0	9.5±6.0
	STASN [36]	6.7±10.6	8.3±8.4	7.5±4.7
	FaceDs [15]	5.1±6.3	6.1±5.1	5.6±5.7
	Ours	5.8±4.9	1.7±2.6	3.7±2.1

generalization capacity over unknown spoof types and our method obtains considerable performance gain over the state-of-the-art. The experimental results demonstrate the superior generalization capacity of the proposed method.

The OULU-NPU dataset is a high quality dataset simulating realistic mobile authentication scenarios. It consists of 4950 videos recorded by different mobile phones with front cameras. OULU-NPU develops four protocols to evaluate the model’s generalization. Table 2 compares the performance of our model with state-of-the-art methods on the OULU-NPU dataset. The state-of-the-art methods we compare with include GRADIANT [5], MILHP [21] Auxiliary [22], STASN [36] and FaceDs [15]. The experimental results show that our model has great generalization ability over unseen environmental conditions, attack mediums, and camera sensors.

4.3 Inter Testing

To demonstrate the generalization capability of our model, we set up inter testing experiments. Specifically, the model is trained on one dataset and then tested on the other dataset. The cross-dataset evaluation is challenging as the data distribution varies a lot on both the live and spoof samples between different datasets. We perform the inter testing on the CASIA-MFSD and Replay-Attack datasets to evaluate the generalization capability of our model. As shown in Table 3, our method achieves the best results under the HTER metric among all the state-of-the-art methods. All the state-of-the-art methods perform the inter test on the CASIA-MFSD and Replay-Attack datasets except that Liu

Table 3. The inter testing results between CASIA-MFSD and Replay-Attack datasets and comparison with the state-of-the-art methods. Results under HTER (%) metric are reported.

Method	Train	Test	Train	Test
	CASIA-MFSD	Replay-Attack	Replay-Attack	CASIA-MFSD
Spectral Cubes [29]	34.4%		50.0%	
CNN [35]	48.5%		45.5%	
LBP [6]	47.0%		39.6%	
Color Texture [7]	30.3%		37.7%	
STASN [36]	31.5%		30.9%	
FaceDs [15]	28.5%		41.1%	
Auxiliary [22]	27.6%		28.4%	
Ours	27.4%		23.7%	

Table 4. The inter testing results from SiW to OULU-NPU dataset and comparison with the state-of-the-art methods. Results under ACER (%) metric are reported.

Protocol	Method	ACER (%)
1	Auxiliary [22]	10.0
	Ours	6.5
2	Auxiliary [22]	14.1
	Ours	11.4
3	Auxiliary [22]	13.8±5.7
	Ours	9.9±8.2
4	Auxiliary [22]	10.0±8.8
	Ours	6.7±6.6

et al. [22] additionally performs an inter test from SiW to OULU. According to Table 3 and Table 4, the proposed method achieves better performance. To conclude, our method outperforms the state-of-the-art methods on all the conducted inter tests. During the inter test, we notice a slight performance drop from the high-resolution dataset (CASIA-MFSD) to the low-resolution dataset (Replay-Attack). As demonstrated in Section 4.5, with high-resolution images as input, the proposed method takes advantage of the abundant texture information, which may be missing on the low-resolution images. On the contrary, the learned spoof cues from low-resolution images can generalize well to high-resolution ones, which is worthy of future research.

4.4 Ablation Studies

Influence of Input Image Resolution. In our model, we sample patches from images as the network input during training. These patches keep the input image resolution and retain as much texture information as possible. In Table 5, we compare resizing the cropped faces to 224×224 with the bilinear interpolation and patching 224×224 the patches from the face images. The experimental result

Table 5. The influence of input image resolution, regression loss (RL), triplet loss (TL) and the auxiliary classification (AC). Results on protocol 1 of OULU-NPU dataset are reported.

Resized input	Patched input	RL on live and spoof	RL on live	TL on E5-D4	TL on E5-SC	AC on spoof cue	AC on overlayed	APCER (%)	BPCER (%)	ACER (%)
✓			✓	✓			✓	2.1	3.3	2.7
	✓		✓				✓	1.7	3.3	2.5
	✓		✓		✓		✓	0.2	4.2	2.3
	✓		✓	✓				2.5	0.8	1.7
	✓	✓		✓			✓	0.8	2.5	1.7
	✓		✓	✓		✓		1.7	0.0	0.8
	✓		✓	✓			✓	0.8	0.0	0.4

shows that patched input strategy achieves better ACER metric. The result is in line with our expectations as our method excavates spoof cues from a single-frame RGB image without extra supervision such as depth or temporal information, the local image information is critical to the performance. On the high resolution dataset like OULU-NPU, the resize operation will blur the image and bring considerable loss in local details, resulting in a performance drop.

Influence of Implicit Supervisions. In the proposed method we adopt the triplet loss (TL) and the auxiliary classification (AC) in a joint supervision manner. In Table 5, we remove the triplet learning and the auxiliary classifier separately and report their performances. It shows that the model is most likely to overfitting if only guided by the binary classification supervision. The joint supervision boosts the performance and results in the minimum ACER. To further validate the influence of the triplet loss on the spoof cue map (SC), we set up a comparison experiment which alternatively imposes the triplet loss on feature maps from E5 to SC. The experimental result shows that imposing triplet loss on feature maps from E5 to D4 achieves better performance.

Influence of Regression Loss. In order to relieve the spoof diversity’s influence on the decision boundary modeling, we set explicit regression loss (RL) only for the live samples while setting no regression loss for the spoof ones. To validate the effectiveness of such design, we set up a comparison experiment which imposes regression loss on both live and spoof samples. To be specific, given an RGB image I as input, the spoof cue map C is a zero map for a live sample while an all-one map for a spoof samples. As shown in Table 5, the proposed design achieves better performance. From our perspective, the predefined explicit regression loss on the spoof samples may violate the intrinsic spoof diversity and makes the network easy to overfitting.

Advantage of Auxiliary Supervision. In Table 5, we evaluate the auxiliary classifier and the residual learning approach’s influence on the network’s performance. The residual learning framework achieves the best ACER performance among all these experiments. For a more intuitive understanding, we present the spoof cue maps from the models in Fig. 3. The auxiliary classifier ampli-

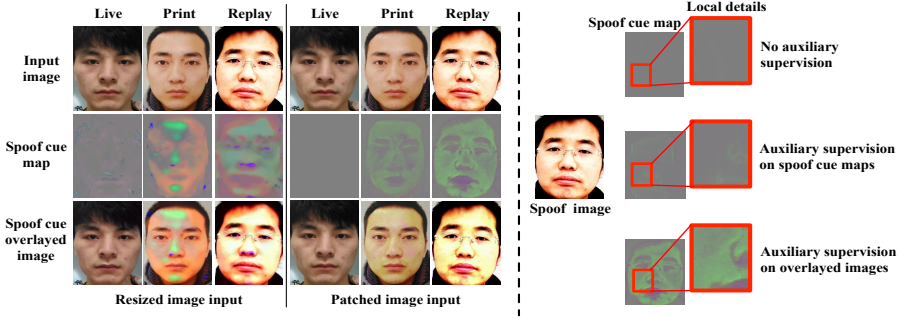


Fig. 3. The illustration of different types of input and the auxiliary supervision’s amplification effect on the spoof cue maps. **Left:** Results with resized and patched images as input during training. **Right:** The spoof cue maps and their local details corresponding to different auxiliary supervision strategies.

fies the spoof cues by enlarging the magnitude of the non-zero values. In our perspective, the correlations between the spoof cue maps and the input images are learned and amplified through the residual learning manner. From the local details column in Fig. 3, our method excavates more local spoof details than other auxiliary supervision manners and obtains more discriminative spoof cue representations.

4.5 Visualizations

Visualization of the Spoof Cues. For a more intuitive understanding of what the spoof cues are, we present the generated spoof cue map in Fig. 3. The spoof cues respond to the differences between the live and spoof samples, which could possibly be the color distortion, moire pattern or other spoof patterns. Fig. 3 shows that the spoof cues exhibit some global pattern when the input image contains a complete face, whereas, for patched input images, the spoof cues contain more local information. In our perspective, the phenomenon is caused by our modeling architecture. Multi-scale skip connections from encoder to decoder and the residual 1×1 Conv in Decoder Residual Block retain image semantics in decoder feature maps well.

Visualization of the Feature Embeddings. We assume the spoof samples are outliers of the live set and develop the joint supervision method to promote the live-live intra-class compactness and the live-spoof inter-class separability. In Fig. 4, we present the feature embeddings of layer D4 using t-SNE [23] on the test set of OULU-NPU protocol 1. It shows that all the live samples fall into a cluster while the spoof samples distribute far from this cluster by a considerable margin, which is in line with our desired distribution of the live closed-set and the spoof open-set.

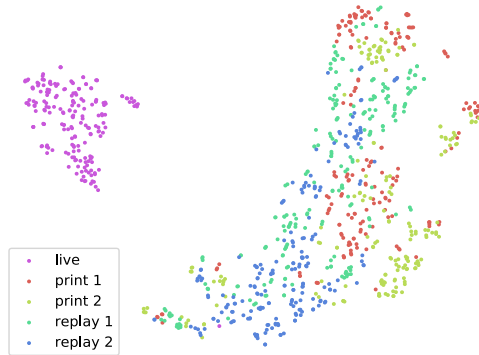


Fig. 4. The t-SNE representation of feature embeddings of layer D4. The samples for demonstration are chosen from the test set of OULU-NPU on protocol 1.

5 Conclusions

We reformulate the FAS in an anomaly detection perspective and assume the live samples belong to a closed-set while the spoof samples are outliers from this closed-set belonging to an open-set. We define the spoof cues as the discriminative live-spoof differences and propose a residual learning framework to learn the generalized spoof cues. Our network consists of a spoof cue generator and an auxiliary classifier. In spoof cue generator we impose weak implicit constraint on the spoof samples to guide the proposed network to learn more generalized spoof cues, the auxiliary classifier makes the spoof cues more discriminative and further boosts the spoof cues’ generalizability. We conduct extensive experiments on popular datasets and the experimental results show that our method achieves state-of-the-art performance for face anti-spoofing.

References

1. Adewumi, A.O., Akinyelu, A.A.: A survey of machine-learning and nature-inspired based credit card fraud detection techniques. *International Journal of System Assurance Engineering and Management* **8**(2), 937–953 (2017)
2. Arashloo, S.R., Kittler, J.: Client-specific anomaly detection for face presentation attack detection. *arXiv preprint arXiv:1807.00848* (2018)
3. Arashloo, S.R., Kittler, J., Christmas, W.: An anomaly detection approach to face spoofing detection: A new formulation and evaluation protocol. *IEEE Access* **5**, 13868–13882 (2017)
4. Au, C.E., Skaff, S., Clark, J.J.: Anomaly detection for video surveillance applications. In: *18th International Conference on Pattern Recognition (ICPR’06)*. vol. 4, pp. 888–891. IEEE (2006)
5. 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). pp. 688–696. IEEE (2017)
6. Boulkenafet, Z., Komulainen, J., Hadid, A.: Face anti-spoofing based on color texture analysis. In: 2015 IEEE international conference on image processing (ICIP). pp. 2636–2640. IEEE (2015)
 7. Boulkenafet, Z., Komulainen, J., Hadid, A.: Face antispoofing using speeded-up robust features and fisher vector encoding. *IEEE Signal Processing Letters* **24**(2), 141–145 (2016)
 8. Boulkenafet, Z., Komulainen, J., Li, L., Feng, X., Hadid, A.: 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). pp. 612–618. IEEE (2017)
 9. Chingovska, I., Anjos, A., Marcel, S.: 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). pp. 1–7. IEEE (2012)
 10. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
 11. Feng, L., Po, L.M., Li, Y., Xu, X., Yuan, F., Cheung, T.C.H., Cheung, K.W.: Integration of image quality and motion cues for face anti-spoofing: A neural network approach. *Journal of Visual Communication and Image Representation* **38**, 451–460 (2016)
 12. 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. pp. 770–778 (2016)
 13. He, X., Zhou, Y., Zhou, Z., Bai, S., Bai, X.: Triplet-center loss for multi-view 3d object retrieval. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1945–1954 (2018)
 14. Hermans, A., Beyer, L., Leibe, B.: In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737* (2017)
 15. Jourabloo, A., Liu, Y., Liu, X.: Face de-spoofing: Anti-spoofing via noise modeling. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 290–306 (2018)
 16. King, D.E.: Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research* **10**(Jul), 1755–1758 (2009)
 17. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
 18. Kollreider, K., Fronthaler, H., Faraj, M.I., Bigun, J.: Real-time face detection and motion analysis with application in liveness assessment. *IEEE Transactions on Information Forensics and Security* **2**(3), 548–558 (2007)
 19. Komulainen, J., Hadid, A., Pietikäinen, M.: Context based face anti-spoofing. In: 2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS). pp. 1–8. IEEE (2013)
 20. Kwon, D., Natarajan, K., Suh, S.C., Kim, H., Kim, J.: An empirical study on network anomaly detection using convolutional neural networks. In: 2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS). pp. 1595–1598. IEEE (2018)
 21. Lin, C., Liao, Z., Zhou, P., Hu, J., Ni, B.: Live face verification with multiple instantialized local homographic parameterization. In: IJCAI. pp. 814–820 (2018)

22. 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*. pp. 389–398 (2018)
23. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(Nov), 2579–2605 (2008)
24. Määttä, J., Hadid, A., Pietikäinen, M.: Face spoofing detection from single images using micro-texture analysis. In: *2011 international joint conference on Biometrics (IJCB)*. pp. 1–7. IEEE (2011)
25. Nikisins, O., Mohammadi, A., Anjos, A., Marcel, S.: On effectiveness of anomaly detection approaches against unseen presentation attacks in face anti-spoofing. In: *2018 International Conference on Biometrics (ICB)*. pp. 75–81. IEEE (2018)
26. Pan, G., Sun, L., Wu, Z., Lao, S.: Eyeblick-based anti-spoofing in face recognition from a generic webcam. In: *2007 IEEE 11th International Conference on Computer Vision*. pp. 1–8. IEEE (2007)
27. Patel, K., Han, H., Jain, A.K.: Secure face unlock: Spoof detection on smartphones. *IEEE transactions on information forensics and security* **11**(10), 2268–2283 (2016)
28. Pérez-Cabo, D., Jiménez-Cabello, D., Costa-Pazo, A., López-Sastre, R.J.: Deep anomaly detection for generalized face anti-spoofing. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops* (2019)
29. Pinto, A., Pedrini, H., Schwartz, W.R., Rocha, A.: Face spoofing detection through visual codebooks of spectral temporal cubes. *IEEE Transactions on image processing* **24**(12), 4726–4740 (2015)
30. 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. *arXiv preprint arXiv:1904.12490* (2019)
31. Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Siddiqui, S.A., Binder, A., Müller, E., Kloft, M.: Deep one-class classification. In: *International conference on machine learning*. pp. 4393–4402 (2018)
32. Shao, R., Lan, X., Yuen, P.C.: Regularized fine-grained meta face anti-spoofing. *arXiv preprint arXiv:1911.10771* (2019)
33. Wang, K., Zhao, Y., Xiong, Q., Fan, M., Sun, G., Ma, L., Liu, T.: Research on healthy anomaly detection model based on deep learning from multiple time-series physiological signals. *Scientific Programming* **2016** (2016)
34. Xu, Z., Li, S., Deng, W.: Learning temporal features using lstm-cnn architecture for face anti-spoofing. In: *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*. pp. 141–145. IEEE (2015)
35. Yang, J., Lei, Z., Li, S.Z.: Learn convolutional neural network for face anti-spoofing. *arXiv preprint arXiv:1408.5601* (2014)
36. 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: *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
37. Zhang, Z., Yan, J., Liu, S., Lei, Z., Yi, D., Li, S.Z.: A face antispoofing database with diverse attacks. In: *2012 5th IAPR international conference on Biometrics (ICB)*. pp. 26–31. IEEE (2012)
38. Zhao, C., Qin, Y., Wang, Z., Fu, T., Shi, H.: Meta anti-spoofing: Learning to learn in face anti-spoofing. *arXiv preprint arXiv:1904.12490* (2019)