# LIVENESS SCORE-BASED REGRESSION NEURAL NETWORKS FOR FACE ANTI-SPOOFING

**Youngjun Kwak** [1 2]  **Minyoung Jung** [3]  **Hunjae Yoo** [1]  **JinHo Shin** [1]  **Changick Kim** [2]

## Abstract

Previous anti-spoofing methods have used either pseudo maps or user-defined labels, and the performance of each approach depends on the accuracy of the third party networks generating pseudo maps and the way in which the users define the labels. In this study, we propose a liveness score-based regression network for overcoming the dependency on third party networks and users. First, we introduce a new labeling technique, called pseudo-discretized label encoding for generating discretized labels indicating the amount of information related to real images. Secondly, we suggest the expected liveness score based on a regression network for training the difference between the proposed supervision and the expected liveness score. Finally, extensive experiments were conducted on four face anti-spoofing benchmarks to verify our proposed method on both intra-and cross-dataset tests. The experimental results denote our approach outperforms previous methods.

## 1. Introduction

Face anti-spoofing (FAS) systems have been successfully established in face authentication, and widely used in online banking, electronic payments, and securities as a crucial technique. Despite its substantial success, FAS still shows vulnerability to various presentation attacks (PAs) such as printed materials, replay-videos, and 3D-masks. To alleviate such vulnerability, previous deep learning-based FAS methods (Liu et al., 2018; Yu et al., 2020b) learn discriminative features for distinguishing real faces against PAs, and such methods mostly treat the FAS problem as a binary classification of whether a given face is real or a spoof, as shown in Fig. 1(a). However, such binary classification-based approaches suffer from non-trivial attacks because they are

[1]KakaoBank Corp., South Korea [2]Department of Electrical Engineering, KAIST, South Korea [3]AIRC, Korea Electronics Technology Institute, South Korea. Correspondence to: Changick Kim <changick@kaist.ac.kr>.
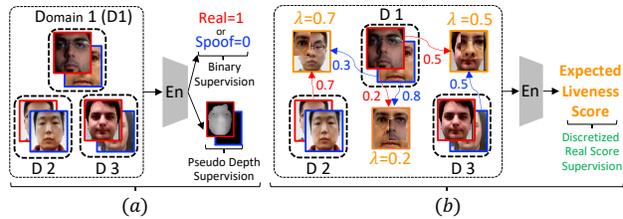
*Figure 1.* Comparison between previous methods and our method for face anti-spoofing. (a) Previous methods utilize either binary supervision to detect spoof cues, or pseudo depth supervision, or both. (b) Our method discretizes binary labels and exchanges real and spoof images for our expected liveness score. The discretized label $\lambda$ indicates the ratio of a real image over an image.

prone to an over-fitting to the training data, resulting in poor generalization (Liu et al., 2018). To mitigate the over-fitting problem, regression-based methods (Feng et al., 2018; Bian et al., 2022; Yu et al., 2021; 2020b) have been proposed, which find sparse evidence for known spoof types and generalize to unseen attacks. For regression-based neural networks, two approaches are considered: First, pseudo-define based supervision (Jiang & Sun, 2022; Bian et al., 2022; Yu et al., 2021; 2020b; Liu et al., 2018; Fang et al., 2021) is designed for context-agnostic discrimination describing the local cues from the pixel level, such as the depth and reflection. For example, a pseudo-map based CNN (Feng et al., 2018) utilizes pseudo-depth supervision using the mean square error (MSE) to reconstruct a sparse depth-map and a flattened-map for a real and spoof image, respectively, as illustrated in Fig. 1(a). And PAL-RW (Fang et al., 2021) is the first approach that utilizes partial pixel-wise labels with face masks to train FAS models. Secondly, user-define based supervision (Jiang Fangling; Wang et al., 2022a) is designed for constrained learning using the relative distances among real and PAs to improve the generalization ability. For instance, ordinal regression (Jiang Fangling) introduces user-defined ordinal labels. Based on the user-defined labels, the model is trained to finely constrain the relative distances among the features of different spoof categories within the latent space. Another example is PatchNet (Wang et al., 2022a), which subdivides binary labels (a real or a spoof) into fine-grained labels (reals or spoofs). Despite previous efforts, we found that the pseudo-define based
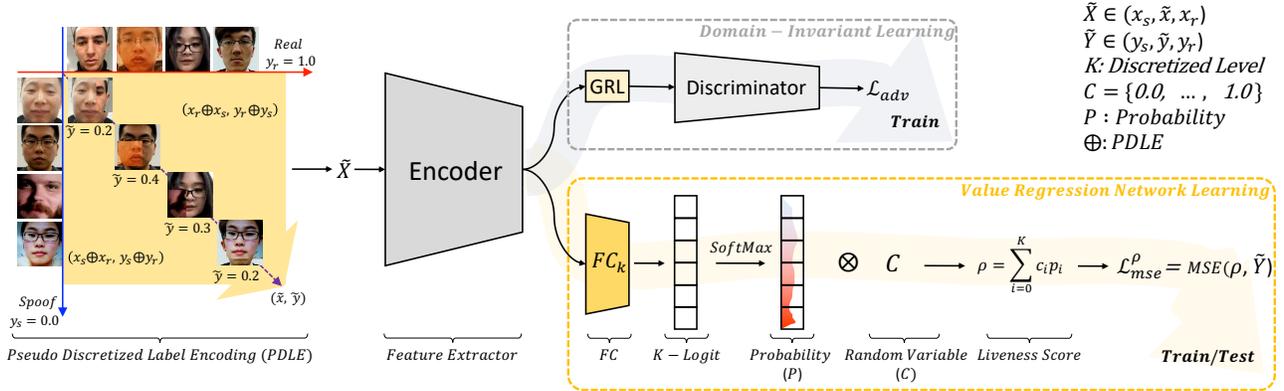
*Figure 2.* Overview of our approach for a value regression neural network. Our framework consists of a label encoding (PDLE) for the data and label expansion, a encoder network for the feature extractor, an expected liveness score estimator for the regression network learning, and a discriminator for the domain-invariant feature learning.

supervisions depend on the accuracy of additional works (e.g., depth- (Feng et al., 2018) and texture-based (Zhang et al., 2018)), and that user-define based supervision relies on that user-specified guides, and the correctness is not guaranteed. In this paper, as described in Fig. 1(b), we introduce a discretized label encoding for increasing data distribution and generating data relationships, which has no dependencies on the prior works. For our proposed label encoding, we present a novel pre-processing method, called the pseudo-discretized label encoding (PDLE) scheme, in which an image is randomly selected in a mini-batch, then the opposite labeled image is also arbitrarily chosen from the whole batch, and then parts of the images are exchanged to generate a new image and its discretized dense label.

Our contributions are as follows:

- We re-formulate face anti-spoofing as a value regression problem that directly optimizes a deep neural network with mean square error for improving performance, instead of using binary cross-entropy.

- We propose a simple yet effective pseudo-discretized label encoding (PDLE), which enforce the regression network to represent the ratio of information of the real image to that of the given input image for a prediction of the liveness score.

- We conduct extensive experiments, and obtain the state-of-the-art and outstanding performance on the intra- and cross-dataset respectively.

## 2. Proposed Method

### 2.1. Overview

For an expansion of the training image and label distribution without an information corruption, we introduce a discretized label encoding schema to preserve the spoof and real cues in the images, and indicate the amount of a real image information over that of the input image. To leverage the PDLE, we propose learning a value regression neural network using the MES between the expected liveness scores and the pseudo-labels. In addition, we apply a domain-invariant learning scheme (GRL) (Ganin & Lempitsky, 2015) as an adversarial training to our regression neural network using the domain labels. The framework of our method is illustrated in Fig. 2.

### 2.2. Pseudo-Discretized Label Encoding

We assume that $X = \{x_s, x_r\} \in \mathbb{R}^{H \times W \times 3}$ and $Y = \{y_s = 0.0, y_r = 1.0\}$ denote the spoof and real color image space and the class label space in each. To sample the discretized labels between $y_s$ and $y_r$, we use the following formula:

$$u \sim \mathcal{U}\{1, K\}; \quad \lambda = \frac{u}{K}, \tag{1}$$

where $u$ is sampled from the discrete uniform distribution $(1, K)$, and $K$ is a pre-defined discretized level, a cardinality of an encoding label $\tilde{Y}$ set, and a number of outputs for the last $FC$ in Fig. 2. $\lambda$ implies a pseudo-discretized label presenting the amount of a partial real image over a whole image. Inspired by CutMix (Yun et al., 2019), we first exchange a real image and a spoof image through a random rectangular box as follows:

$$\tilde{x} = M \odot x_a + (1 - M) \odot x_b, \text{ where } y_a \neq y_b$$

$$\tilde{y} = \begin{cases} 1 - \lambda, & \text{if } x_a = x_r \\ \lambda, & \text{otherwise,} \end{cases} \tag{2}$$

where $M \in \{0, 1\}^{H \times W}$ is a random rectangular mask based on $\lambda$, with 0 and 1 indicating inside and outside the mask. $\odot$ is an element-wise multiplication operator, and $x_a$ is an anchor to choose a sample from a mini-batch, whereas $x_b$ is the opposite sample selected from the entire training set. $\tilde{x}$ indicates the exchanged image, and $\tilde{y}$ is the pseudo-discretized label determined based on whether $x_a$ is a real image or not. We exchange between images with different labels ($y_a \neq y_b$) to expand data and label distribution. As shown in Fig. 2, we use $\tilde{X} \in (x_s, \tilde{x}, x_r)$ and $\tilde{Y} \in (y_s, \tilde{y}, y_r)$

as the training data and the supervision for the regression network to learn the liveness score.

## 2.3. Expected Liveness Score

Let $\mathbb{P} : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^K$ denote the probability of the liveness evidence estimated using $SoftMax$, $FC_k$, and $Encoder(\tilde{X})$, as illustrated in Fig. 2. We employ $K$ in Eq. 1 to formulate a random variable $C$ with a finite list $\{c_0, ..., c_K\}$ whose the $i^{th}$ element $c_i$ is denoted as follows:

$$c_i = \begin{cases} 0.0, & \text{if } i = 0 \\ interval \times i, & \text{if } i > 0 \text{ and } i < K \\ 1.0, & \text{if } i = K \end{cases} \quad (3)$$

where $interval = \lceil \frac{y_r - y_s}{K} \rceil$. The random variable $c_i$ and its probability $p_i$ are exploited to calculate the expected liveness score as follows:

$$\mathbb{E}[C] = \rho = \sum_{i=0}^{K} c_i * p_i, \quad (4)$$

where $p_i$ is the $i^{th}$ element of $P$ which is the predicted probability vector of real cues from the input $\tilde{X}$. We write $\mathbb{E}[C]$ with $\rho$, which is calculated using the sum over the element-wise multiplication between the random variables and their corresponding probabilities.

## 2.4. Objective Function

Our objective function is defined as follows:

$$L^{\rho}_{mse} = -\frac{1}{N} \sum_{j=1}^{N} (\rho_j - \tilde{Y}_j)^2, \quad (5)$$

where $N$ is a mini-batch size, and $\tilde{Y}_j$ and $\rho_j$ are the $j^{th}$ supervision and expected liveness score in the mini-batch. We calculate the distance between $\tilde{Y}_j$ and $\rho_j$ for our main objective function $L^{\rho}_{mse}$.

To further improve the performance, we exploit not only a regression network but also an adversarial learning technique GRL (Ganin & Lempitsky, 2015). Finally, our overall loss function can be formulated as follows:

$$L_{final} = \alpha * L^{\rho}_{mse} + (1 - \alpha) * L_{adv}, \quad (6)$$

where $L^{\rho}_{mse}$ is a liveness score-based regression training loss and $L_{adv}$ is an adversarial training loss for jointly learning our livensss score-based regression neural network. $\alpha$ is a non-negative parameter to balance the importance of two losses, and we empirically set $\alpha$ to 0.5.

## 3. Experiments

We demonstrate the effectiveness of the proposed approach on an intra- and cross-dataset. Based on the experimental results, the characteristics of our algorithm will be discussed in this section.

## 3.1. Datasets and Metrics

**Datasets.** We employed four public datasets, OULU-NPU (labeled O) (Boulkenafet et al., 2017), CASIA-FASD (la-

beled C) (Zhang et al., 2012), Replay-Attack (labeled I) (Chingovska et al., 2012), and MSU-MFSD (labeled M) (Wen et al., 2015) for our experiments. OULU-NPU is a high-resolution database with four protocols for validating the improved performance on the intra-dataset. The videos of each dataset are recorded under different scenarios with various cameras and subjects, and they are used for cross-dataset testing to validate the generalization ability for testing data with unconstrained distribution shifts.

**Evaluation Metrics.** We utilized average classification error rate (ACER) for the intra-dataset testing on OULU-NPU. The half total error rate (HTER) and area under curve (AUC) are measured for the cross-dataset testing protocols.

## 3.2. Implementation Details

**Primitive Data Preparation and Augmentation.** Because the four FAS datasets are in video format, we extracted images at certain intervals. After obtaining the images, we used RetinaFace (Deng et al., 2019) to detect faces, and then cropped and resized the color image to a resolution of $256 \times 256$. Data augmentation, including horizontal flipping and random cropping, was used for training, and center cropping was employed for testing. And we empirically set $K$ to 10 for our approach after testing variant $K$ as depicted in Fig. 3.

**Experimental Setting.** To train the FAS task, we used ResNet18 (He et al., 2015) as the encoder with the Adam optimizer under an initial learning rate and weight decay of 1e-4 and 2e-4, respectively, for all testing protocols. We trained the models with a batch size of 32 and a max epoch of 200, whereas decreasing the learning rate through an exponential LR with a gamma of 0.99. For the domain labels on the intra-dataset, we used the number of sessions in each protocol.

## 3.3. Intra-Dataset Testing on OULU-NPU

OULU-NPU has four protocols for evaluating the generalization ability under mobile scenarios with previously unseen sensors and spoof types. As shown in Table. 1, our PDLE approach presents the best performance for all protocols, and the expected liveness scores clearly validate the ability to generalize better latent embedding features. In particular, our proposed PDLE achieves the significant performance im-

*Table 1.* Evaluation results for ACER (%) in comparison with the previous methods and the proposed **PDLE** approach within the intra-dataset (OULU-NPU protocols).

| Method | Protocol 1 ACER(%) | Protocol 2 ACER(%) | Protocol 3 ACER(%) | Protocol 4 ACER(%) |
|---|---|---|---|---|
| Auxiliary (Liu et al., 2018) | 1.6 | 2.7 | 2.9±1.5 | 9.5±6.0 |
| CDCN (Yu et al., 2020b) | 1.0 | 1.45 | 2.3±1.4 | 6.9±2.9 |
| FaceDs (Amin Jourabloo*, 2018) | 1.5 | 4.3 | 3.6±1.6 | 5.6±5.7 |
| DC-CDN (Yu et al., 2021) | 0.4 | 1.3 | 1.9±1.1 | 4.3±3.1 |
| LMFD-PAD (Fang et al., 2022) | 1.5 | 2.0 | 3.4±3.1 | 3.3±3.1 |
| NAS-FAS (Yu et al., 2020a) | 0.2 | **1.2** | 1.7±0.6 | 2.9±2.8 |
| PatchNet (Wang et al., 2022a) | **0** | **1.2** | 1.18±1.26 | 2.9±3.0 |
| Ours | **0** | **1.2** | **0.96±1.03** | **0.63±1.04** |

*Table 2.* Comparison results of cross-domain testing on MSU-MFSD (M), CASIA-MFSD (C), Replay-Attack (I), and OULU-NPU (O). PE and LE mean patch-exchange and label-encoding, respectively. **Bold** and *italic* denote the best results among Res-18 and Res-50 based methods in each.

| Method | Network | O&C&I to M | | O&M&I to C | | O&C&M to I | | I&C&M to O | |
|---|---|---|---|---|---|---|---|---|---|
| | | HTER(%) | AUC(%) | HTER(%) | AUC(%) | HTER(%) | AUC(%) | HTER(%) | AUC(%) |
| NAS-FAS (Yu et al., 2020a) | NAS | 19.53 | 88.63 | 16.54 | 90.18 | 14.51 | 93.84 | 13.80 | 93.43 |
| NAS-FAS w/ D-Meta (Yu et al., 2020a) | NAS | 16.85 | 90.42 | 15.21 | 92.64 | 11.63 | 96.98 | 13.16 | 94.18 |
| DRDG (George & Marcel, 2021) | DenseNet | 12.43 | 95.81 | 19.05 | 88.79 | 15.56 | 91.79 | 15.63 | 91.75 |
| ANRL (Liu et al., 2021) | - | 10.83 | 96.75 | 17.83 | 89.26 | 16.03 | 91.04 | 15.67 | 91.90 |
| LMFD-PAD (Fang et al., 2022) | Res-50 | 10.48 | 94.55 | 12.50 | 94.17 | 18.49 | 84.72 | 12.41 | *94.95* |
| DBEL (Jiang & Sun, 2022) | Res-50 | 8.57 | 95.01 | 20.26 | 85.80 | *13.52* | *93.22* | 20.22 | 88.48 |
| HFN+MP (Cai et al., 2022) | Res-50 | *5.24* | 97.28 | *9.11* | 96.09 | 15.35 | 90.67 | *12.04* | 94.26 |
| CAFD (Huang et al., 2022) | Res-18 | 11.64 | 95.27 | 17.51 | 89.98 | 15.08 | 91.92 | 14.27 | 93.04 |
| SSDG-R (Jia et al., 2020) | Res-18 | 7.38 | 97.17 | 10.44 | 95.94 | 11.71 | 96.59 | 15.61 | 91.54 |
| SSAN-R (Wang et al., 2022b) | Res-18 | 6.67 | 98.75 | **10.00** | **96.67** | 8.88 | 96.79 | 13.72 | 93.63 |
| PatchNet (Wang et al., 2022a) | Res-18 | 7.10 | 98.46 | 11.33 | 94.58 | 13.40 | 95.67 | 11.82 | 95.07 |
| Ours w/o PE&LE | Res-18 | 10.83 | 94.58 | 15.08 | 91.14 | 14.50 | 93.55 | 13.88 | 93.16 |
| Ours w/o PE | Res-18 | 10.41 | 94.93 | 13.59 | 91.04 | 11.17 | 93.92 | 12.50 | 94.35 |
| Ours w/o LE | Res-18 | 9.58 | 94.47 | 12.47 | 92.28 | 12.25 | 94.55 | 13.29 | 93.62 |
| Ours | Res-18 | **5.41** | **98.85** | 10.05 | 94.27 | **8.62** | **97.60** | **11.42** | **95.52** |

provement for protocol 4 (unseen lighting, spoof type, and sensor type). The results demonstrate that the effectiveness to train a liveness score-based regression neural network using the amount of swapping as pseudo-discrete labels. Note that our proposed PDLE improves the overall ACER performance over the previous SOTA (PatchNet (Wang et al., 2022a)) approach.

### 3.4. Cross-Dataset Testing

To evaluate our proposed method, we select three out of four datasets to train and use the remaining one for testing, denoted by $\{\cdot\&\cdot\&\cdot\}$ to $\{\bullet\}$. We compare our proposed method with the latest methods as shown in Table 2. Among ResNet-18 based methods, we found that our method shows outstanding performance on testing the O&C&I to M, O&C&M to I, and I&C&M to O protocols, and the other protocol O&M&I to C displays the very competitive performance. When comparing to the ResNet-50 based method HFN+MP (Cai et al., 2022), our approach shows competitive performance on testing datasets O&C&I to M and O&M&I to C which contain a variety of image resolutions, and superior performance on testing O&C&M to I and I&C&M to O whose images are collected from various capture devices unlike other datasets. By split testing on each capture device in the dataset C, we found that our method show relatively the low performance on low quality images (93.73% AUC) compared to normal (94.79% AUC) and high quality (96.47% AUC) images. This result proves that the proposed method achieves satisfactory performance on all protocols because our liveness score-based regression network estimates probabilities of the real cues under various presentation attacks.

### 3.5. Ablation Study

We conducted ablation studies on cross-dataset testing to explore the contribution of each component in our method, as depicted in Table 2. To analyze the effect of discretization, we separated the proposed PDLE into patch exchange

(PE) and label encoding (LE). And we confirmed that each of them is the essential element for improving performance, and also observed the best performance when both were used. In addition, we verified the influence of the pre-defined $K$ in PDLE for determining the representation power of the liveness against an input image. As shown in Fig. 3, we tested various values of $K$ on the O&C&M to I protocol to investigate the impact of $K$ on AUC. With $K$ between 2 and 17, our method outperforms the baseline.
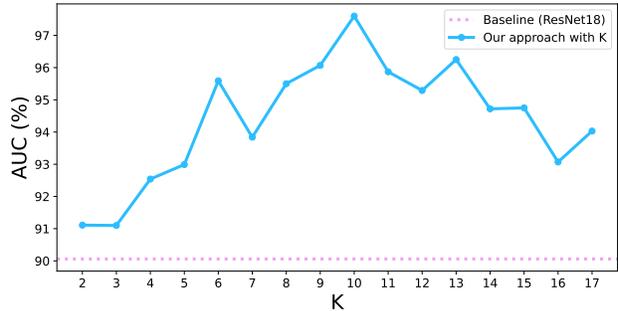


*Figure 3.* Ablation study on the discretized level $K$.

## 4. Conclusion

In this paper, we have proposed the PDLE approach for training a face anti-spoofing regression model. The regression model allows the probability to estimate our liveness score. Our approach not only has the effect of a data augmentation because different labels and domains are densely exchanged, new data combinations are also created, which results in the improved domain generalization. Through our experiments, we confirm the effectiveness, robustness, and generalization of the proposed PDLE and expected liveness score.

## 5. Acknowledgements

# References

Amin Jourabloo*, Yaojie Liu*, X. L. Face de-spoofing: Anti-spoofing via noise modeling. In ECCV 2018, Munich, Germany, 2018.

Bian, Y., Zhang, P., Wang, J., Wang, C., and Pu, S. Learning multiple explainable and generalizable cues for face anti-spoofing. In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2310–2314, 2022. doi: 10.1109/ICASSP43922.2022.9747677.

Boulkenafet, Z., Komulainen, J., Li, L., Feng, X., and Hadid, A. OULU-NPU: A mobile face presentation attack database with real-world variations. In FG 2017, pp. 612–618, 2017. doi: 10.1109/FG.2017.77.

Cai, R., Li, Z., Wan, R., Li, H., Hu, Y., and Kot, A. C. Learning meta pattern for face anti-spoofing. IEEE TIFS, 2022.

Chingovska, I., Anjos, A., and Marcel, S. On the effectiveness of local binary patterns in face anti-spoofing. In IOSIG, pp. 1–7, 2012.

Deng, J., Guo, J., Yuxiang, Z., Yu, J., Kotsia, I., and Zafeiriou, S. Retinaface: Single-stage dense face localisation in the wild. In arxiv, 2019.

Fang, M., Boutros, F., Kuijper, A., and Damer, N. Partial attack supervision and regional weighted inference for masked face presentation attack detection. IEEE FG, 2021.

Fang, M., Damer, N., Kirchbuchner, F., and Kuijper, A. Learnable multi-level frequency decomposition and hierarchical attention mechanism for generalized face presentation attack detection. In WACV, pp. 1131–1140. IEEE, 2022.

Feng, Y., Wu, F., Shao, X., Wang, Y., and Zhou, X. Joint 3d face reconstruction and dense alignment with position map regression network. In ECCV, 2018.

Ganin, Y. and Lempitsky, V. Unsupervised domain adaptation by backpropagation, 2015.

George, A. and Marcel, S. Cross modal focal loss for rgbd face anti-spoofing. In CVPR, 2021.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. CoRR, abs/1512.03385, 2015. URL http://arxiv.org/abs/1512.03385.

Huang, H., Xiang, Y., Yang, G., Lv, L., Li, X., Weng, Z., and Fu, Y. Generalized face anti-spoofing via cross-adversarial disentanglement with mixing augmentation. In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2939–2943, 2022. doi: 10.1109/ICASSP43922.2022.9746716.

Jia, Y., Zhang, J., Shan, S., and Chen, X. Single-side domain generalization for face anti-spoofing. In CVPR, 2020.

Jiang, J. and Sun, Y. Depth-based ensemble learning network for face anti-spoofing. In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2954–2958, 2022. doi: 10.1109/ICASSP43922.2022.9747840.

Jiang Fangling, Liu Pengcheng, Z. X.-D. Ordinal regression with representative feature strengthening for face anti-spoofing.

Liu, S., Zhang, K., Yao, T., Bi, M., Ding, S., Li, J., Huang, F., and Ma, L. Adaptive normalized representation learning for generalizable face anti-spoofing. ACM, 2021.

Liu, Y., Jourabloo, A., and Liu, X. Learning deep models for face anti-spoofing: Binary or auxiliary supervision. In CVPR, June 2018.

Wang, C.-Y., Lu, Y.-D., Yang, S.-T., and Lai, S.-H. Patchnet: A simple face anti-spoofing framework via fine-grained patch recognition. In CVPR, pp. 20281–20290, June 2022a.

Wang, Z., Wang, Z., Yu, Z., Deng, W., Li, J., Gao, T., and Wang, Z. Domain generalization via shuffled style assembly for face anti-spoofing. CoRR, abs/2203.05340, 2022b. doi: 10.48550/arXiv.2203.05340.

Wen, D., Han, H., and Jain, A. K. Face spoof detection with image distortion analysis. IEEE Transactions on Information Forensics and Security, 10(4):746–761, 2015. doi: 10.1109/TIFS.2015.2400395.

Yu, Z., Wan, J., Qin, Y., Li, X., Li, S. Z., and Zhao, G. Nasfas: Static-dynamic central difference network search for face anti-spoofing. In TPAMI, 2020a.

Yu, Z., Zhao, C., Wang, Z., Qin, Y., Su, Z., Li, X., Zhou, F., and Zhao, G. Searching central difference convolutional networks for face anti-spoofing. In CVPR, 2020b.

Yu, Z., Qin, Y., Zhao, H., Li, X., and Zhao, G. Dual-cross central difference network for face anti-spoofing. In IJCAI, 2021.

Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. Cutmix: Regularization strategy to train strong classifiers with localizable features. In ICCV, 2019.

Zhang, X., Ng, R., and Chen, Q. Single image reflection separation with perceptual losses. In IEEE Conference on Computer Vision and Pattern Recognition, 2018.

Zhang, Z., Yan, J., Liu, S., Lei, Z., Yi, D., and Li, S. Z. A
   face antispoofing database with diverse attacks. In ICB,
   pp. 26–31, 2012. doi: 10.1109/ICB.2012.6199754.