

Multi-Modal Face Anti-Spoofing Based on Central Difference Networks

Zitong Yu¹, Yunxiao Qin², Xiaobai Li¹, Zezheng Wang, Chenxu Zhao³, Zhen Lei⁴, Guoying Zhao^{1*}

¹CMVS, University of Oulu ²Northwestern Polytechnical University

³Mininglamp Academy of Sciences, Mininglamp Technology ⁴Authenmetric

{zitong.yu, xiaobai.li, guoying.zhao}@oulu.fi, {qyxqyx}@mail.nwpu.edu.cn

{zhaochenxu}@mininglamp.com, {zhen.lei}@authenmetric.com

Abstract

Face anti-spoofing (FAS) plays a vital role in securing face recognition systems from presentation attacks. Existing multi-modal FAS methods rely on stacked vanilla convolutions, which is weak in describing detailed intrinsic information from modalities and easily being ineffective when the domain shifts (e.g., cross attack and cross ethnicity). In this paper, we extend the central difference convolutional networks (CDCN) [39] to a multi-modal version, intending to capture intrinsic spoofing patterns among three modalities (RGB, depth and infrared). Meanwhile, we also give an elaborate study about single-modal based CDCN. Our approach won the first place in Track Multi-Modal as well as the second place in Track Single-Modal (RGB) of ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020 [20]. Our final submission obtains $1.02 \pm 0.59\%$ and $4.84 \pm 1.79\%$ ACER in Track Multi-Modal and Track Single-Modal (RGB), respectively. The codes are available at <https://github.com/ZitongYu/CDCN>.

1. Introduction

Face recognition has been widely used in many interactive artificial intelligence systems for its convenience (e.g., access control, face payment and device unlock). However, vulnerability to presentation attacks (PAs) curtails its reliable deployment. Merely presenting printed images or videos to the biometric sensor could fool face recognition systems. Typical examples of presentation attacks are print, video replay, and 3D masks. For the reliable use of face recognition systems, face anti-spoofing (FAS) methods are important to detect such presentation attacks.

In recent years, several hand-crafted feature based [3, 4, 7, 15, 28, 27] and deep learning based [38, 33, 29, 22, 12, 34, 2, 8, 9] methods have been proposed for presentation attack detection (PAD). On one hand, the classical hand-

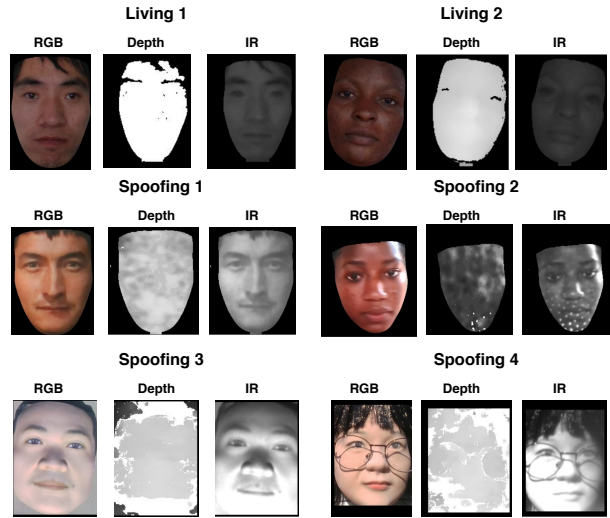


Figure 1. Examples of living and spoofing faces from CASIA-SURF CeFA dataset [21].

crafted descriptors (e.g., local binary pattern (LBP) [3]) leverage local relationship among the neighbours as the discriminative features, which is robust for describing the detailed invariant information (e.g., color texture, moiré pattern and noise artifacts) between the living and spoofing faces. On the other hand, due to the stacked convolution operations with nonlinear activation, the convolutional neural networks (CNN) hold strong representation abilities to distinguish the bona fide from PAs. However, CNN based methods focus on the deeper semantic features, which are weak in describing detailed intrinsic information between living and spoofing faces and easily being ineffective when acquisition conditions varies (e.g., light illumination and camera type). In order to solve this issue, central difference convolutional networks (CDCN) is developed [39] for single-modal (RGB) FAS task and achieves state-of-the-art performance on several benchmark datasets. Although the state-of-the-art single-modal FAS methods are robust in some existing testing protocols, it is still challenging when

*denotes corresponding author

encountering new kinds of domain shift (e.g., cross ethnicity).

Recently, a large-scale cross-ethnicity face anti-spoofing dataset, the CASIA-SURF CeFA [21], is established, which covers three ethnicities, three modalities, 1607 subjects, and 2D plus 3D attack types. Some typical examples are shown in Fig. 1. The most challenging protocol 4 (simultaneously cross-attack and cross-ethnicity) is utilized for ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020 [20]. The baseline results in CASIA-SURF CeFA dataset [21] indicate: 1) multiple modalities (i.e., RGB, depth and infrared (IR)) fusion is more robust than using an arbitrary single modal, and 2) the multi-modal result, only $31.8 \pm 10.0\%$ ACER in protocol 4, is barely satisfactory. Hence it is necessary to explore more effective multi-modal FAS methods for cross-attack and cross-ethnicity testing.

Motivated by the discussions above, we first analyze how different modality influences the performance of CDCN. Then we extend CDCN to a multi-modal version, intending to capture intrinsic spoofing patterns among various modalities. Our contributions include:

- We are the first to utilize CDCN for depth and infrared modalities based FAS and analyze how CDCN performs with these two modalities. Besides considering CDCN as a single-modal network, we extend it to a multi-modal version, which captures rich discriminative clues among modalities and represents invariant intrinsic patterns across ethnicities and attacks.
- Our approach won the first place in Track Multi-Modal¹ as well as the second place in Track Single-Modal (RGB)² of ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020 [20].

2. Related Work

In this section, we first introduce some recent progress in the single-modal FAS community; and then demonstrate few recent works about multi-modal FAS. Finally, classical convolution operators for vision tasks are presented.

Single-Modal Face Anti-Spoofing. Traditional single-modal face anti-spoofing methods usually extract hand-crafted features from the RGB facial images to capture the spoofing patterns. Several classical local descriptors such as LBP [3, 7], SIFT [27], SURF [5], HOG [15] and DoG [28] are utilized to extract frame level features while video level methods usually capture dynamic clues like dynamic texture [14], micro-motion [32] and eye blinking [24]. More recently, a few deep learning based methods are proposed for both frame level and video level face anti-spoofing. For

frame level methods [39, 29, 16, 26, 9, 12], deep CNN models are utilized to extract features in a binary-classification setting. In contrast, auxiliary depth supervised FAS methods [2, 22] are introduced to learn more detailed information effectively. On the other hand, several video level CNN methods are presented to exploit the dynamic spatio-temporal [33, 34, 19] or rPPG [17, 22, 18, 36, 37, 31] features for PAD. Despite achieving state-of-the-art performance, single-modal methods are easily influenced by unseen domain shift (e.g., cross ethnicity and cross attack types) and not robust for challenging cases (e.g., harsh environment and realistic attacks).

Multi-Modal Face Anti-Spoofing. There are also few works for multi-modal face anti-spoofing. Zhang et al. [40] take ResNet18 as the backbone and propose a three-stream network, where the input of each stream is RGB, Depth and IR face images, respectively. Then, these features are concatenated and passed to the last two residual blocks. Aleksandr et al. [25] also consider the similar fusion network with three streams. ResNet34 is chosen as the backbone and multi-scale features are fused at all residual blocks. Tao et al. [30] present a multi-stream CNN architecture called FaceBagNet. In order to enhance the local detailed representation ability, patch-level images are adopted as inputs. Moreover, modality feature erasing operation is designed to prevent overfitting and obtain more robust modal-fused features. All previous methods just consider standard backbone (ResNet) with stacked vanilla convolutions for multiple modalities, which might be weak in representing the intrinsic features between living and spoofing faces.

Convolution Operators. The convolution operator is commonly used in extracting basic visual features in deep learning framework. Recently extensions to the vanilla convolution operator have been proposed. In one direction, classical local descriptors (e.g., LBP [1] and Gabor filters [11]) are considered into convolution design. Representative works include Local Binary Convolution [13] and Gabor Convolution [23], which are proposed for saving computational cost and enhancing the resistance to the spatial changes, respectively. Recently, Yu et al. propose Central Difference Convolution (CDC) [39], which is suitable for FAS task because of its excellent representation ability for detailed intrinsic patterns. Another direction is to modify the spatial scope for aggregation. Two related works are dilated convolution [35] and deformable convolution [6]. However, these convolution operators are always designed for RGB modality, it is still unknown how they perform for depth and IR modalities.

In order to overcome the above-mentioned drawbacks and fill in the blank, we extend the state-of-the-art single-modal network CDCN to a multi-modal version for challenging cross-ethnicity and cross-attack FAS task.

¹<https://competitions.codalab.org/competitions/23318>

²<https://competitions.codalab.org/competitions/22151>

3. Methodology

In this section, we will first introduce CDC [39] as a preliminary in Section 3.1, then demonstrate our single-modal and multi-modal neural architectures in Section 3.2 and Section 3.3, respectively. At last the supervision signals and loss functions are presented in Section 3.4.

3.1. Preliminary: CDC

The feature maps and convolution can be represented in 3D shape (2D spatial domain and extra channel dimension) in modern deep learning frameworks. For simplicity, all convolutions in this paper are described in 2D while extension to 3D is straightforward.

Vanilla Convolution. There are two main steps in the 2D spatial convolution: 1) *sampling* local receptive field region \mathcal{R} over the input feature map x ; 2) *aggregation* of sampled values via weighted summation. Hence, the output feature map y can be formulated as

$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n), \quad (1)$$

where p_0 denotes current location on both input and output feature maps while p_n enumerates the locations in \mathcal{R} . For instance, local receptive field region for convolution operation with 3×3 kernel and dilation 1 is $\mathcal{R} = \{(-1, -1), (-1, 0), \dots, (0, 1), (1, 1)\}$.

Central Difference Convolution. For FAS task, the discriminative and robust features indicate fine-grained living/spoofing patterns and environment invariant clues, respectively. Local gradient operator (e.g., basic element in local binary pattern (LBP) [3]), as a residual and difference term, is able to capture rich detailed patterns and not easily affected by external changes.

Inspired by LBP [3], we introduce central difference context into vanilla convolution to enhance its representation and generalization capacity. Similar to vanilla convolution, central difference convolution also consists of two steps, i.e., *sampling* and *aggregation*. The sampling step is similar to that in vanilla convolution while the aggregation step is different: central difference convolution prefers to aggregate the center-oriented gradient of sampled values. Thus Eq. (1) becomes

$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot (x(p_0 + p_n) - x(p_0)). \quad (2)$$

When $p_n = (0, 0)$, the gradient value always equals to zero with respect to the central location p_0 itself.

As both the intensity-level semantic information and gradient-level detailed message are crucial for distinguishing the living and spoofing faces, which indicates that combining vanilla convolution with central difference convolution might be a feasible manner to provide more robust

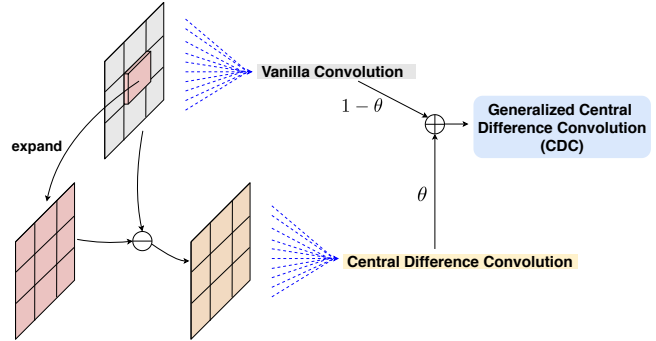


Figure 2. Generalized central difference convolution (CDC).

modeling capacity. As illustrated in Fig. 2, we generalize central difference convolution as

$$y(p_0) = \theta \cdot \underbrace{\sum_{p_n \in \mathcal{R}} w(p_n) \cdot (x(p_0 + p_n) - x(p_0))}_{\text{central difference convolution}} + (1 - \theta) \cdot \underbrace{\sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n)}_{\text{vanilla convolution}}, \quad (3)$$

where hyperparameter $\theta \in [0, 1]$ tradeoffs the contribution between intensity-level and gradient-level information. The higher value of θ means the more importance of central difference gradient information. Similar to [39], we refer to this generalized central difference convolution as **CDC**.

3.2. Single-Modal CDCN

We follow the similar configuration ‘CDCN++’ [39] as our single-modal backbone, including low-mid-high level cells and Multiscale Attention Fusion Module (MAFM). In the consideration of the large-scale training data in CASIA-SURF CeFA dataset, we set the initial channel number as 80 instead of 64. The specific network is shown in Fig. 3(a). Single-modal face image with size $256 \times 256 \times 3$ is taken as the network input and the output is the predicted 32×32 grayscale mask.

3.3. Multi-Modal CDCN

We adopt the configuration ‘CDCN’ [39] as the backbone of each modality branch as we find the MAFM would drop the performance when using multi-modal fusion. As illustrated in Fig. 3(b), the backbone network of each modality branch is not shared. Thus each branch is able to learn modality-aware features independently. The multi-level features from each modality branch are fused via concatenation. Finally, the two head layers aggregate the multi-modal features and predict the grayscale mask.

As the feature-level fusion strategy might not be optimal for all protocols, we also try two other fusion strategies: 1)

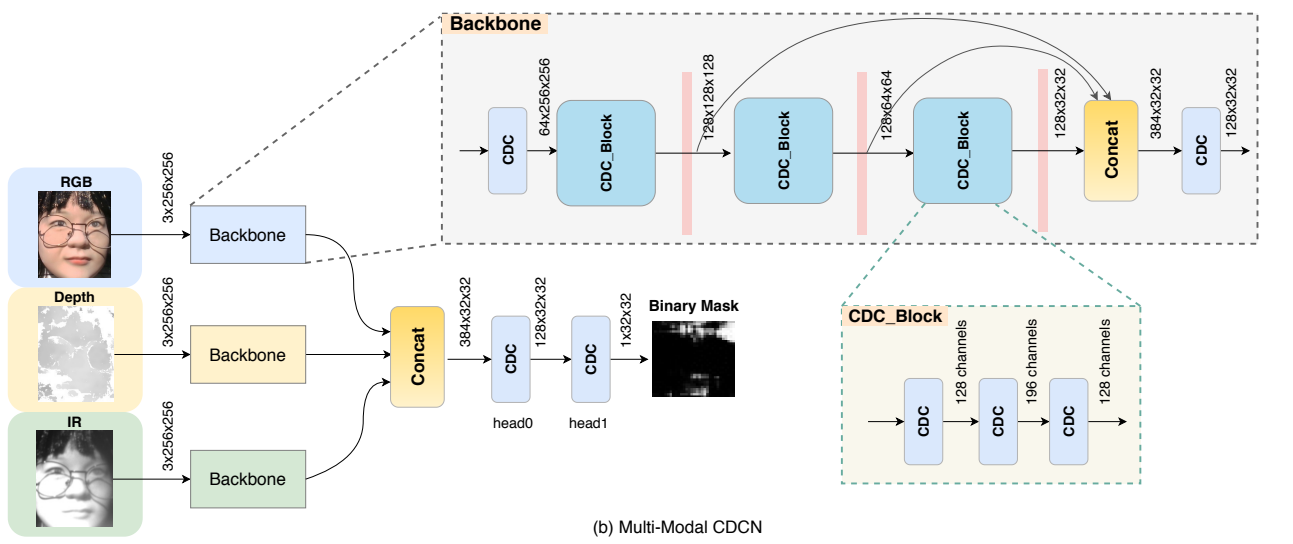
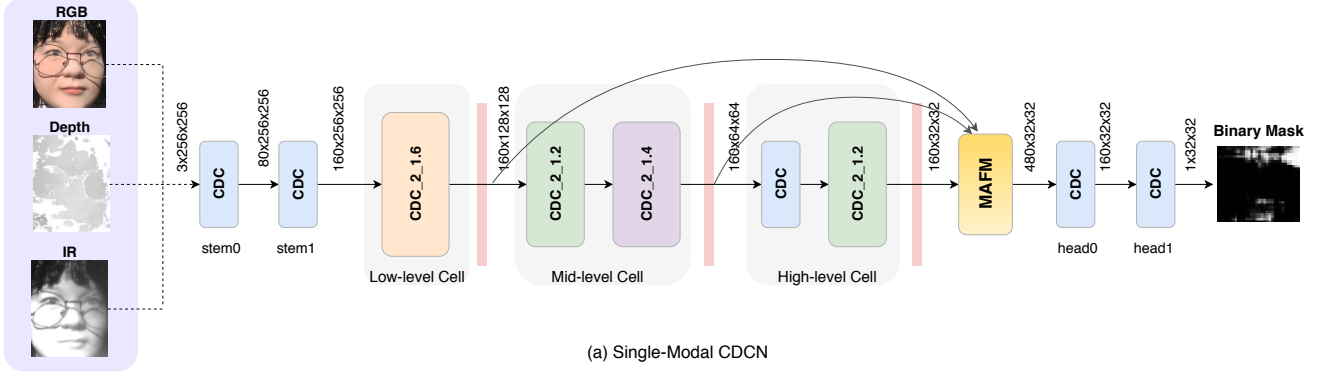


Figure 3. The architecture of (a) single-modal and (b) multi-modal CDCN. The red thin rectangle denotes a max pool layer with stride 2. ‘CDC_2_r’ means using two stacked CDC to increase channel number with ratio r first and then decrease back to the original channel size.

input-level fusion via concatenating three-modal inputs to $256 \times 256 \times 9$ directly, and 2) score-level fusion via weighting the predicted score from each modality. For these two fusion strategies, the architecture of single-modal CDCN (see Fig. 3(a)) is used. The corresponding ablation study will be shown in Section 4.4.

3.4. Supervision

Compared with traditional guidance from the binary scalar score, pixel-wise supervision [9] helps to learn more discriminative patterns between living and spoofing faces. As a result, our network prefers to predict 32×32 grayscale mask instead of traditional scalar score. In terms of ground truth label, we generate the binary mask via simply set the non-zero pixel value to ‘1’ because the intensity values of non-face background have already been ‘0’ in CASIA-SURF CeFA dataset.

For the loss function, mean square error loss \mathcal{L}_{MSE} is

$$\begin{aligned}
 & \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix}_0, \begin{bmatrix} 0 & 1 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix}_1, \begin{bmatrix} 0 & 0 & 1 \\ 0 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix}_2, \begin{bmatrix} 0 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix}_3, \\
 & \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & 0 \end{bmatrix}_4, \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 1 & 0 & 0 \end{bmatrix}_5, \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 1 & 0 \end{bmatrix}_6, \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix}_7
 \end{aligned}$$

Figure 4. The kernel K_n^{CDL} in contrastive depth loss..

utilized for pixel-wise supervision, which is formulated:

$$\mathcal{L}_{MSE} = \frac{1}{H \times W} \sum_{i \in H, j \in W} (B_{pre}(i,j) - B_{gt}(i,j))^2, \quad (4)$$

where H, W denote the height and width of the binary mask, respectively, and B_{pre} and B_{gt} mean the predicted grayscale mask and ground truth binary mask, respectively. Moreover, for the sake of fine-grained supervision needs in FAS task, contrastive depth loss (CDL) \mathcal{L}_{CDL} [33] is considered to help the networks learn more detailed features.

Table 1. Ablation study of the hyperparameter θ with RGB modality.

Single-Modal CDCN	Protocol 4@1			Protocol 4@2			Protocol 4@3			Overall
	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	ACER(%)
$\theta=0.5$	12.61	4.0	8.31	6.67	2.0	4.33	4.56	8.5	6.53	6.39
$\theta=0.6$	11.67	8.0	9.83	10.56	3.0	6.78	3.89	5.0	4.44	7.02
$\theta=0.7$	12.83	1.25	7.04	13.33	2.0	7.67	3.72	3.0	3.36	6.02
$\theta=0.8$	14.33	1.5	7.92	10.0	6.25	8.13	3.83	7.25	5.54	7.19
$\theta=0.9$	11.17	2.5	6.83	21.33	5.75	13.54	3.56	7.5	5.53	8.63

Table 2. Results of Single-Modal CDCN ($\theta=0.7$) with different modalities.

Modality	Protocol 4@1			Protocol 4@2			Protocol 4@3			Overall
	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	ACER(%)
RGB	12.83	1.25	7.04	13.33	2.0	7.67	3.72	3.0	3.36	6.02
Depth	5.22	1.25	3.24	2.72	0.5	1.61	4.94	1.75	3.35	2.73
IR	1.56	1.0	1.28	27.72	0.25	13.99	29.56	0.5	15.03	10.1

Table 3. Best submission result in Track Single-Modal (RGB).

Method	Protocol 4@1			Protocol 4@2			Protocol 4@3			Overall
	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	ACER(%)
SD-Net [21]	-	-	-	-	-	-	-	-	-	35.2±5.8
Ours (Single-Modal)	11.17	2.5	6.83	6.67	2.0	4.33	3.72	3.0	3.36	4.84±1.79

CDL can be formulated as

$$\mathcal{L}_{CDL} = \frac{\sum_{i \in H, j \in W, n \in N} (K_n^{CDL} \odot B_{pre(i,j)} - K_n^{CDL} \odot B_{gt(i,j)})^2}{H \times W \times N}, \quad (5)$$

where K_n^{CDL} is the n -th contrastive convolution kernel, and N denotes the kernel numbers. The details of the kernels ($N = 8$) can be found in Fig. 4. Finally, the overall loss $L_{overall}$ can be formulated as $\mathcal{L}_{overall} = \mathcal{L}_{MSE} + \mathcal{L}_{CDL}$.

4. Experiments

In this section, extensive experiments are performed to demonstrate the effectiveness of our method. In the following, we sequentially describe the employed datasets & metrics (Sec. 4.1), implementation details (Sec. 4.2), results (Sec. 4.3 - 4.4) and visualization (Sec. 4.5).

4.1. Datasets and Metrics

CASIA-SURF CeFA Dataset [21]. CASIA-SURF CeFA aims to provide with the largest up to date face anti-spoofing dataset to allow for the evaluation of the generalization performance cross-ethnicity and cross-attacks. It consists of 2D and 3D attack subsets. For the 2D attack subset, it includes print and video-reply attacks, and three ethnicities (African, East Asian and Central Asian) with two attacks (print face from cloth and video-replay). Each ethnicity has 500 subjects. Each subject has one real sample, two fake samples of print attack captured in indoor and outdoor, and 1 fake sample of video-replay. In total, there are 18000 videos (6000 per modality).

There are four evaluation protocols in CASIA-SURF CeFA for cross-ethnicity, cross-attack, cross-modality, and cross-ethnicity & cross-attack testing. In this paper, our experiments are all conducted on the most challenging protocol 4 (cross-ethnicity & cross-attack), which has been

utilized for ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020.

Performance Metrics. Three metrics, i.e., Attack Presentation Classification Error Rate (APCER), Bona Fide Presentation Classification Error Rate (BPCER), and Average Classification Error Rate (ACER) [10] are utilized for performance comparison. They can be formulated as

$$\begin{aligned} APCER &= \frac{FP}{TN + FP}, \\ BPCER &= \frac{FN}{FN + TP}, \\ ACER &= \frac{APCER + BPCER}{2}, \end{aligned} \quad (6)$$

where FP , FN , TN and TP denote the false positive, false negative, true negative and true positive sample numbers, respectively. ACER is used to determine the final ranking in ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020.

4.2. Implementation Details

Our proposed method is implemented with Pytorch. In the training stage, models are trained with Adam optimizer and the initial learning rate and weight decay are $1e-4$ and $5e-5$, respectively. We train models with 50 epochs while learning rate halves every 20 epochs. The batch size is 8 on a P100 GPU. In the testing stage, we calculate the mean value of the predicted grayscale map as the final score.

4.3. Single-Modal Testing

In this subsection, we give the ablation study about the hyperparameter θ with RGB modality firstly. Then based on the optimal θ for CDCN, we test depth and IR modalities.

Table 4. Ablation study of fusion strategies for multi-modal CDCN. We only report the results tried in the FAS challenge.

Modality	Protocol 4@1			Protocol 4@2			Protocol 4@3		
	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)
Feature-level fusion	0.33	0.5	0.42	5.89	3.25	4.57	4.22	3.25	3.74
Input-level fusion	0.5	3.75	2.13	5.67	1.5	3.58	2.61	3.25	2.93
Score-level fusion	-	-	-	1.39	0.75	1.07	1.44	1.75	1.6

Table 5. Best submission result in Track Multi-Modal.

Method	Protocol 4@1			Protocol 4@2			Protocol 4@3			Overall
	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	APCER(%)	BPCER(%)	ACER(%)	ACER(%)
PSMM-Net [21]	33.3	15.8	24.5	78.2	8.3	43.2	50.0	5.5	27.7	31.8±10.0
Ours (Multi-Modal)	0.33	0.5	0.42	1.39	0.75	1.07	1.44	1.75	1.6	1.02±0.59

Finally, we summarize our best submission results in Track Single Modal (RGB) on ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020.

Impact of θ with RGB modality. As shown in Table 1, the best overall performance (ACER=6.02%) is achieved when $\theta = 0.7$, which is consistent with the evidence in [39]. As for the sub-protocols, $\theta = 0.9$, $\theta = 0.5$ and $\theta = 0.7$ obtain the lowest ACER in protocol 4@1 (6.83%), 4@2 (4.33%) and 4@3 (3.36%), respectively.

Results of Depth and IR modalities. Table 2 shows the results of different modalities using single-modal CDCN when $\theta = 0.7$. It is surprising that the performance varies a lot across modalities. The IR modality performs the best in protocol 4@1 (testing without Africa) but the worst in protocol 4@2 and 4@3 (testing with Africa), indicating that the IR modality generalizes poorly for unseen Africa ethnicity. Compared with RGB and IR modalities, the depth modality is more robust and discriminative in most cases (e.g., print attacks in testing stage) because the 3D depth shape is quite distinguishable between living and print faces. The excellent overall performance indicates central difference convolution is not only suitable for RGB modality, but also for IR and depth modalities.

Best Submission Result in Track Single-Modal (RGB). Our best submission result (4.84±1.79% ACER) is shown in Table 3, which wins the second place in Track Single-Modal (RGB) on ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020. This final result is combined with the best sub-protocols results (i.e., $\theta = 0.9$, 0.5 and 0.7, respectively).

4.4. Multi-Modal Testing

In this subsection, three fusion strategies are studied in multi-modal testing. Then the best submission results in Track Multi-Modal will be presented.

Multi-Modal Fusion Strategies. As shown in Table 4, our proposed multi-modal CDCN (i.e., feature-level fusion with three modalities) achieves the lowest ACER (0.42%) in protocol 4@1. When using the concatenated inputs with three modalities (input-level fusion), the CDCN could obtain comparable performance with the single-modal results in Table 2. However, it still causes the performance drops

compared with the best single-modal results (i.e., IR modality for protocol 4@1, depth modality for protocol 4@2 and protocol 4@3). It also reflects the issue for both feature- and input-level fusion, i.e., simple fusion with concatenation might be sub-optimal because it is weak in representing and selecting the importance of modalities. It is worth exploring more effective fusion methods (e.g., attention mechanism for modalities) in future.

Based on the prior results in Table 2, we weight the results of RGB and depth modalities averagely as the score-level fusion (i.e., $fusion_score = 0.5 * RGB_score + 0.5 * depth_score$). As shown in Table 4 (the third row), this simple ensemble strategy helps to boost the performance significantly. Compared with single-depth modality, score-level fusion gives 0.54% and 1.13% ACER improvements for protocol 4@2 and 4@3, respectively.

Best Submission Result in Track Multi-Modal. Table 3 shows our best submission result (1.02±0.59% ACER), which wins the first place in Track Multi-Modal on ChaLearn FAS Attack Detection Challenge@CVPR2020. This final result is combined with the best sub-protocols results (i.e., feature-level fusion for protocol 4@1 while score-level fusion for protocol 4@2 and 4@3).

4.5. Feature Visualization

The visualizations of CDCN with three modalities are shown in Fig. 5. On one hand, it is clear that the low-level, mid-level and high-level features in CDCN are distinguishable between living and spoofing faces among all three modalities. In terms of low-level features, the living have more detailed texture (especially in IR modality). As for the high-level features, the living face regions are purer and plainer while the spoofing ones are with more spoofing/noise patterns.

On the other hand, depth and IR modalities are complementary to RGB modality and helpful for robust liveness detection. We can see from the last row in Fig. 5 that CDCN fails to detect spoofing1 only using RGB input while spoofing1 could be accurately detected by depth or IR inputs.

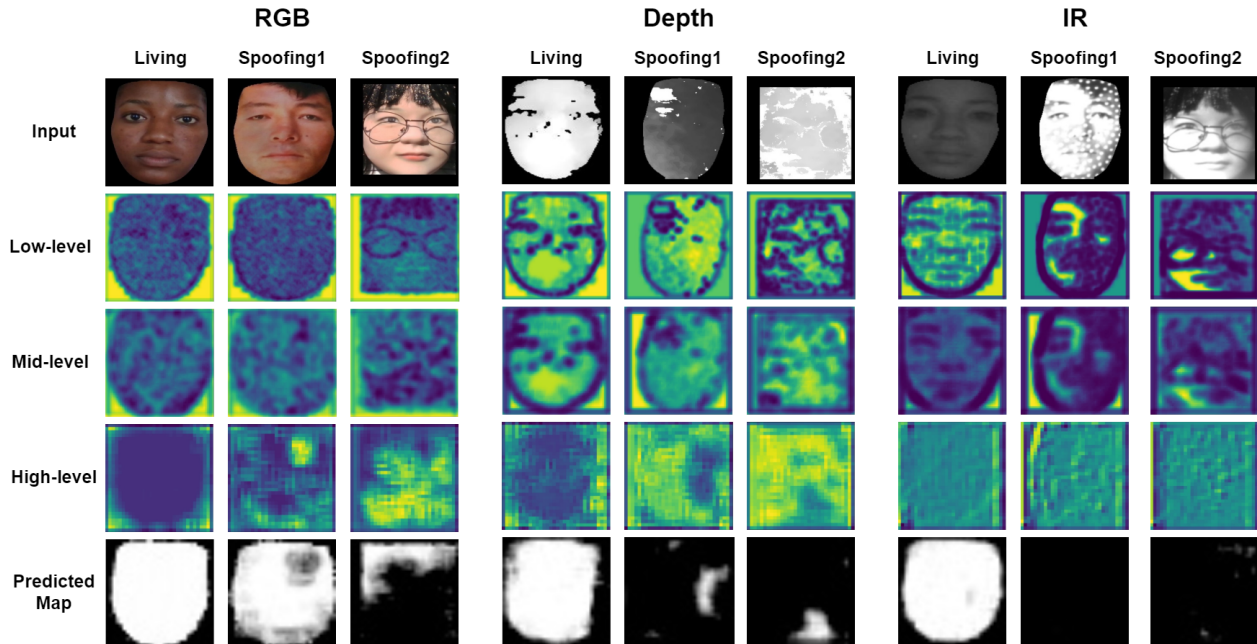


Figure 5. Visualization of CDCN with three modalities.

5. Conclusion

In this paper, we give an elaborate study about the applications of central difference convolutional networks (CDCN) [39] for multiple modalities in face anti-spoofing (FAS) task. The experimental results indicate the effectiveness of CDCN for both single-modal and multi-modal FAS. The proposed approach wins the first place in Track Multi-Modal as well as the second place in Track Single-Modal (RGB) of ChaLearn Face Anti-spoofing Attack Detection Challenge@CVPR2020.

6. Acknowledgement

This work was supported by the Academy of Finland for project MiGA (Grant 316765), ICT 2023 project (Grant 328115), Infotech Oulu and the Chinese National Natural Science Foundation Projects (Grant No. 61876178). As well, the authors acknowledge CSCIT Center for Science, Finland, for computational resources.

References

- [1] Timo Ahonen, Abdenour Hadid, and Matti Pietikainen. Face description with local binary patterns: Application to face recognition. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (12):2037–2041, 2006.
- [2] Yousef Atoum, Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. Face anti-spoofing using patch and depth-based cnns. In *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pages 319–328, 2017.
- [3] Zinelabidine Boulkenafet, Jukka Komulainen, and Abdenour Hadid. Face anti-spoofing based on color texture analysis. In *IEEE international conference on image processing (ICIP)*, pages 2636–2640, 2015.
- [4] 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.
- [5] Zinelabidine Boulkenafet, Jukka Komulainen, and Abdenour Hadid. Face antispoofing using speeded-up robust features and fisher vector encoding. *IEEE Signal Processing Letters*, 24(2):141–145, 2017.
- [6] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017.
- [7] Tiago de Freitas Pereira, André Anjos, José Mario De Martino, and Sébastien Marcel. Lbp- top based countermeasure against face spoofing attacks. In *Asian Conference on Computer Vision*, pages 121–132, 2012.
- [8] Junying Gan, Shanlu Li, Yikui Zhai, and Chengyun Liu. 3d convolutional neural network based on face anti-spoofing. In *ICMIP*, pages 1–5, 2017.
- [9] Anjith George and Sébastien Marcel. Deep pixel-wise binary supervision for face presentation attack detection. In *International Conference on Biometrics*, number CONF, 2019.
- [10] international organization for standardization. Iso/iec jtc 1/sc 37 biometrics: Information technology biometric presentation attack detection part 1: Framework. In <https://www.iso.org/obp/ui/iso>, 2016.
- [11] Anil K Jain and Farshid Farrokhnia. Unsupervised texture segmentation using gabor filters. *Pattern recognition*, 24(12):1167–1186, 1991.

- [12] Amin Jourabloo, Yaojie Liu, and Xiaoming Liu. Face de-spooing: Anti-spoofing via noise modeling. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 290–306, 2018.
- [13] Felix Juefei-Xu, Vishnu Naresh Boddeti, and Marios Savvides. Local binary convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 19–28, 2017.
- [14] Jukka Komulainen, Abdenour Hadid, and Matti Pietikäinen. Face spoofing detection using dynamic texture. In *Asian Conference on Computer Vision*, pages 146–157. Springer, 2012.
- [15] Jukka Komulainen, Abdenour Hadid, and Matti Pietikainen. Context based face anti-spoofing. In *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*, pages 1–8, 2013.
- [16] Lei Li, Xiaoyi Feng, Zinelabidine Boulkenafet, Zhaoqiang Xia, Mingming Li, and Abdenour Hadid. An original face anti-spoofing approach using partial convolutional neural network. In *IPTA*, pages 1–6, 2016.
- [17] Xiaobai Li, Jukka Komulainen, Guoying Zhao, Pong-Chi Yuen, and Matti Pietikäinen. Generalized face anti-spoofing by detecting pulse from face videos. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 4244–4249. IEEE, 2016.
- [18] Bofan Lin, Xiaobai Li, Zitong Yu, and Guoying Zhao. Face liveness detection by rppg features and contextual patch-based cnn. In *Proceedings of the 2019 3rd International Conference on Biometric Engineering and Applications*, pages 61–68. ACM, 2019.
- [19] Chen Lin, Zhouyingcheng Liao, Peng Zhou, Jianguo Hu, and Bingbing Ni. Live face verification with multiple instantiated local homographic parameterization. In *IJCAI*, pages 814–820, 2018.
- [20] Ajian Liu, Xuan Li, Sergio Escalera, Escalante Wan, Sergio Hugo Jair, Meysam Madadi, Yuan Zhuo, Xiaogong Yu, Zichang Tan, Qi Yuan, Ruikun Yang, Benjia Zhou, Guodong Guo, and Stan Z Li. Cross-ethnicity face anti-spoofing recognition challenge: A review. *arXiv preprint arXiv:2003.05136*, 2020.
- [21] Ajian Liu, Zichang Tan, Xuan Li, Jun Wan, Sergio Escalera, Guodong Guo, and Stan Z Li. Casia-surf cefa: A benchmark for multi-modal cross-ethnicity face anti-spoofing. *arXiv preprint arXiv:2003.05136*, 2020.
- [22] Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. Learning deep models for face anti-spoofing: Binary or auxiliary supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 389–398, 2018.
- [23] Shangzhen Luan, Chen Chen, Baochang Zhang, Jungong Han, and Jianzhuang Liu. Gabor convolutional networks. *IEEE Transactions on Image Processing*, 27(9):4357–4366, 2018.
- [24] Gang Pan, Lin Sun, Zhaohui Wu, and Shihong Lao. Eyeblink-based anti-spoofing in face recognition from a generic webcam. In *IEEE International Conference on Computer Vision*, pages 1–8, 2007.
- [25] Aleksandr Parkin and Oleg Grinchuk. Recognizing multi-modal face spoofing with face recognition networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [26] Keyurkumar Patel, Hu Han, and Anil K Jain. Cross-database face antispoofing with robust feature representation. In *Chinese Conference on Biometric Recognition*, pages 611–619, 2016.
- [27] Keyurkumar Patel, Hu Han, and Anil K Jain. Secure face unlock: Spoof detection on smartphones. *IEEE transactions on information forensics and security*, 11(10):2268–2283, 2016.
- [28] Bruno Peixoto, Carolina Michelassi, and Anderson Rocha. Face liveness detection under bad illumination conditions. In *ICIP*, pages 3557–3560. IEEE, 2011.
- [29] Yunxiao Qin, Chenxu Zhao, Xiangyu Zhu, Zezheng Wang, Zitong Yu, Tianyu Fu, Feng Zhou, Jingping Shi, and Zhen Lei. Learning meta model for zero-and few-shot face anti-spoofing. *AAAI*, 2020.
- [30] Tao Shen, Yuyu Huang, and Zhijun Tong. Facebagnet: Bag-of-local-features model for multi-modal face anti-spoofing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [31] Jingang Shi, Iman Alikhani, Xiaobai Li, Zitong Yu, Tapio Seppänen, and Guoying Zhao. Atrial fibrillation detection from face videos by fusing subtle variations. *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [32] Talha Ahmad Siddiqui, Samarth Bharadwaj, Tejas I Dhamecha, Akshay Agarwal, Mayank Vatsa, Richa Singh, and Nalini Ratha. Face anti-spoofing with multifeature videolet aggregation. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 1035–1040. IEEE, 2016.
- [33] Zezheng Wang, Zitong Yu, Chenxu Zhao, Xiangyu Zhu, Yunxiao Qin, Qiusheng Zhou, Feng Zhou, and Zhen Lei. Deep spatial gradient and temporal depth learning for face anti-spoofing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [34] Xiao Yang, Wenhan Luo, Linchao Bao, Yuan Gao, Dihong Gong, Shibao Zheng, Zhifeng Li, and Wei Liu. Face anti-spoofing: Model matters, so does data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [35] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [36] Zitong Yu, Xiaobai Li, and Guoying Zhao. Remote photoplethysmograph signal measurement from facial videos using spatio-temporal networks. In *Proc. BMVC*, pages 1–12, 2019.
- [37] Zitong Yu, Wei Peng, Xiaobai Li, Xiaopeng Hong, and Guoying Zhao. Remote heart rate measurement from highly compressed facial videos: an end-to-end deep learning solution with video enhancement. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 151–160, 2019.
- [38] Zitong Yu, Yunxiao Qin, Xiangqing Xu, Chenxu Zhao, Zezheng Wang, Zhen Lei, and Guoying Zhao. Auto-

fas: Searching lightweight networks for face anti-spoofing. *ICASSP*, 2020.

- [39] Zitong Yu, Chenxu Zhao, Zezheng Wang, Yunxiao Qin, Zhuo Su, Xiaobai Li, Feng Zhou, and Guoying Zhao. Searching central difference convolutional networks for face anti-spoofing. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [40] Shifeng Zhang, Xiaobo Wang, Ajian Liu, Chenxu Zhao, Jun Wan, Sergio Escalera, Hailin Shi, Zezheng Wang, and Stan Z Li. A dataset and benchmark for large-scale multi-modal face anti-spoofing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 919–928, 2019.