

DRL-FAS: A Novel Framework Based on Deep Reinforcement Learning for Face Anti-Spoofing

Rizhao Cai, Haoliang Li, Shiqi Wang, Changsheng Chen, and Alex C. Kot

Abstract—Inspired by the philosophy employed by human beings to determine whether a presented face example is genuine or not, i.e., to glance at the example globally first and then carefully observe the local regions to gain more discriminative information, for the face anti-spoofing problem, we propose a novel framework based on the Convolutional Neural Network (CNN) and the Recurrent Neural Network (RNN). In particular, we model the behavior of exploring face-spoofing-related information from image sub-patches by leveraging deep reinforcement learning. We further introduce a recurrent mechanism to learn representations of local information sequentially from the explored sub-patches with an RNN. Finally, for the classification purpose, we fuse the local information with the global one, which can be learned from the original input image through a CNN. Moreover, we conduct extensive experiments, including ablation study and visualization analysis, to evaluate our proposed framework on various public databases. The experiment results show that our method can generally achieve state-of-the-art performance among all scenarios, demonstrating its effectiveness.

Index Terms—Face anti-spoofing, deep learning, reinforcement learning

I. INTRODUCTION

FACE recognition techniques have been increasingly deployed in everyday scenarios for authentication purposes, such as mobile devices unlocking and door control access. Compared with other biometric information, using faces for authentication is more user-friendly as face verification is non-intrusive, and the face images can be feasibly captured with mobile phone cameras. However, it has been widely recognized that state-of-the-art face recognition systems are still vulnerable to spoofing attacks. Attackers can easily hack a face recognition system by presenting a spoofing face of a client to the system’s camera, where a spoofing face could be a face mask and a face image shown by a printed photo or by a digital display. Therefore, reliable Face Anti-Spoofing (FAS) techniques are highly desired and essential for developing secure face recognition systems.

The past few years have witnessed much progress in the FAS problem. Traditionally, in either the Spatial or Fourier space, various techniques have been proposed to extract hand-crafted features with image descriptors as representations [2], [3], [2], [4], [5], [6]. These features are usually used to train a Support Vector Machine (SVM) to classify genuine or spoofing examples. However, these features are insufficiently discriminative because those descriptors (e.g., Local Binary Pattern) are not originally designed for the FAS problem.

Recently, deep-learning-based methods, which aim to learn discriminative representations in an end-to-end manner, have shown evidence to be more effective in countermeasures

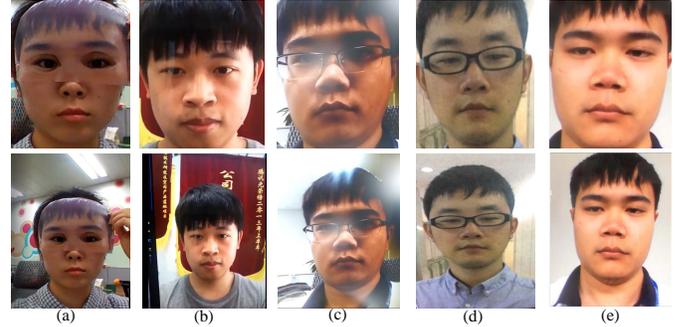


Fig. 1. Presentation attack examples in different scenarios. The examples are from the ROSE-YOUTU Face Liveness Detection Database [1]. The spoofing face examples are shown at the bottom of each column, and the cropped faces are shown at the top. (a) an example of paper mask which is not close enough to the acquisition camera such that the paper boundary can be spotted. (b) an example of video attack which is not close enough to the acquisition camera such that the bezel of the presentation medium (in the left) can be seen. (c) an example of video attack which is acquired with a camera close enough. Although no bezel is visible, obvious spoofing clues (e.g. reflections) appear around the face. (d) an example of print attack. (e) an example of replay attack. The examples in (d)&(e) are launched carefully such that no paper boundary, bezel and obvious spoofing clues can be seen.

against spoofing attacks than the traditional methods. Yang *et al.* [7] firstly introduce the Convolutional Neural Network (CNN) for the FAS task. They train an AlexNet-based model [8], extract features from the model’s last layer and learn an SVM with binary labels (“genuine” or “spoofing”) for classification. Besides using binary labels, Liu *et al.* [9] seek for the auxiliary supervision signals. They use auxiliary techniques to extract pseudo depth maps and remote PhotoPlethysmoGraphy (rPPG) signals from the RGB images for supervision to boost the training. It is also reported that the Recurrent Neural Network (RNN) can be used to utilize temporal information from sequential frames for face anti-spoofing [10], [11]. However, one limitation of the aforementioned techniques is that the learned feature representations may overfit to the properties of a particular database. For example, depth information can benefit face anti-spoofing when the suspicious input is in 2D format (e.g., printed photo, screen display), but it is likely to fail to counter mask attacks, which are with 3D information (e.g., Fig. 1(a)). To learn more spoofing-discriminative representations and alleviate the overfitting effect, we propose a novel two-branch framework based on CNN and RNN.

A. Motivation

The motivations behind this work are inspired by 1) the observation that spoofing clues can appear in various ways and 2) how human beings can act to predict whether a presented face

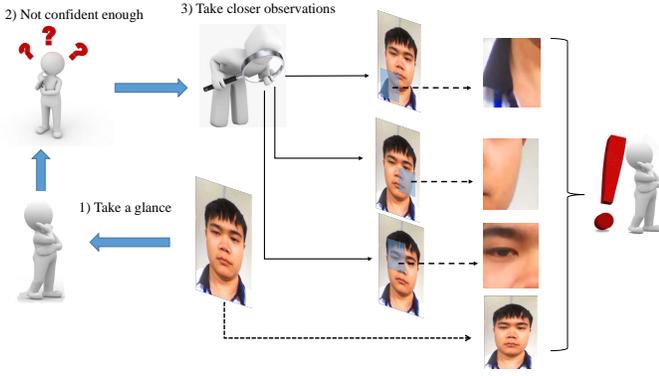


Fig. 2. Illustration of how a testee assesses the liveness of a face example without any apparent distortion. 1) Firstly, the testee will take a glance at the given face example. 2) The testee may not tell that it is a spoofing face with certainty because it looks similar to a genuine face. 3) To confirm the assessment, the testee will take closer observations to carefully search subtle spoofing clues. After several steps of observations, the testee can provide a more accurate assessment by increasingly discovering spoofing clues. Therefore, having observed the example globally and locally, the testee can make a more reliable assessment than a glance without any closer observation.

example is genuine or spoofing. Regarding the first motivation, we show motivating examples in Fig. 1, which indicates that spoofing clues can be diverse. In occasional cases, spoofing clues are visually salient, such as the boundaries of a printed photo, the bezel of a digital display [12], and reflections, which are respectively shown in Fig. 1(a), Fig. 1(b) and Fig. 1(c). These clues can be easily spotted and used by human beings to assess the examples as “spoofing” even without further careful observation. However, in most cases, human beings may give prediction with less certainty as the aforementioned clues may be inconspicuous if an attacker carefully launches the spoofing attack. For instance, no paper boundary, bezel, or reflection appears in Fig. 1(d) and Fig. 1(e). Moreover, the visual quality of Fig. 1(d) and Fig. 1(e) is better than that of Fig. 1(a), Fig. 1(b) and Fig. 1(c). In other words, Fig. 1(d) and Fig. 1(e) look more similar to genuine faces, and thus human beings may not tell the difference with only a glance. Thus, to mine more information to confirm their assessment, human beings would carefully delve into local sub-patches to explore fine-scale and subtle spoofing clues. Fig. 2 portrays such behavior, which reveals our second motivation.

Under these two motivations, we propose a two-branch framework, DRL-FAS, that jointly exploits global and local features based on CNN, RNN, and deep reinforcement learning (DRL), for the face anti-spoofing (FAS) problem. Fig. 3 elaborates on this framework, which corresponds to Fig. 2. Firstly, we treat human beings’ glance at an example as the procedure of extracting global features. As such, we train a CNN to learn global information through the entire frames from video data. Then, we treat the following closer observations at suspicious sub-patches as the procedure of extracting local features. To model such observation behavior, we leverage reinforcement learning to learn a policy model that predicts locations of suspicious sub-patches and learn local information there with RNN. Finally, since human beings can benefit from both global and local information for better prediction, the extracted global and local features are fused

for classification.

The contributions of this work are three-fold:

- We propose a novel framework based on CNN and RNN for the FAS problem. Our framework aims to extract and fuse the global and local features. While many of the previous works used RNN to leverage temporal information from video frames, we take advantage of RNN to memory information from all “observations” from sub-patches to reinforce extracted local features gradually.
- To explore spoofing-specific local information, we leverage the advantage of reinforcement learning to discover suspicious areas where discriminative local features can be extracted. To the best of our knowledge, in the field of FAS, this is the first attempt to introduce reinforcement learning for the optimization.
- We conduct extensive experiments using six benchmark databases to evaluate our method. As shown in Section IV-C, our method can perform better than the schemes that either use global or local features. Moreover, our proposed method can generally achieve state-of-the-art performance compared with other methods.

II. RELATED WORKS

A. Traditional Face Anti-Spoofing

Most of the traditional FAS techniques focus on designing handcrafted features and learning classifiers with learning methods such as SVM. Texture-based methods are based on the assumption that there are differences in texture between genuine faces and spoofing faces due to the inherent nature of different materials. In the Fourier spectrum, Tan *et al.* [5] propose to use Difference-of-Gaussian (DoG) features to describe the frequency disturbance caused by the recapture. Besides, Gragnaniello *et al.* [6] propose to use Local Phase Quantization (LPQ) to analyze texture distortion through the phase of images. Also, texture descriptors, such as Local Binary Pattern (LBP), Scale-Invariant Feature Transform (SIFT), are used in the spatial domain to extract features to represent such disparities [13], [2], [14], [15], [4]. In addition, due to the distortion caused during the recapture process, spoofing examples usually have lower visual quality compared with genuine ones. Motivated by this observation, the FAS community also has proposed to detect spoofing attack examples by assessing the input image quality [16], [17], [18]. Apart from analyzing a single image, methods based on sequential video frames are also proposed to utilize the information from the temporal space, such as eye blinking, lip moves [19], [20], [21], and motion blurring effect [22]. Although such methods can be effective against photo attacks, they cannot counter video replay attacks where such movement information can exist in a given video.

B. Deep-Learning-Based Face Anti-Spoofing

Recently, deep learning has dominated the computer vision community, so as the FAS field. Yang *et al.* [7] are the first to apply CNN to the FAS problem. The authors train a deep model based on AlexNet [8] architecture and extract features from the model’s last layer to train an SVM classifier.

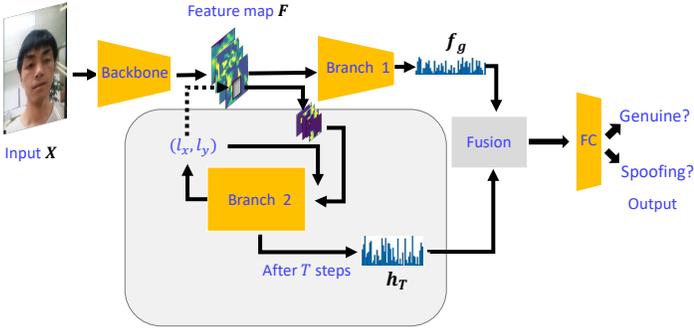


Fig. 3. Overview of the proposed framework. The backbone and Branch 1 networks are based on CNN, and Branch 2 is based on RNN. The working pipeline is as follows: 1) The backbone network processes the input X into the feature map F . 2) Subsequently, Branch 1 extracts global features from the entire F . Meanwhile, Branch 2 extracts local features from sub-patches of F recursively. 3) Finally, the extracted global and local features are fused for the classification in the last fully-connected (FC) layer.

However, this is just a straightforward application of AlexNet, and the improvement is limited compared with the handcrafted features. After that, more CNN-based methods are proposed [23], [24], [25]. Moreover, RNN (e.g., Long Short-Term Memory networks [26], Gated Recurrent Unit [27]) can also be used for the FAS problem by leveraging temporal information from sequential video frames [28], [10]. By far, the aforementioned methods merely use binary labels (“genuine” and “spoofing”) for training. Other than that, the methods in [29], [9], [30] utilize extra information for auxiliary supervision, such as depth. For example, Atoum *et al.* [29] introduce auxiliary depth information for spoofing detection. They hold the idea that 2D spoofing examples of printed papers and screens are flat and thus lack 3D information. As such, they train a CNN-based depth estimator, by which the output toward genuine and spoofing faces are 3D depth maps and flat maps, respectively. Although the performance is shown to be improved with depth as auxiliary supervision, such depth-based methods may not work efficiently when a paper mask attack with depth information is launched.

C. Cross-Domain Face Anti-Spoofing

The variety of capturing settings, such as different cameras, environment illuminations, presentation mediums, etc., can lead to the domain shift problem [1]. To be specific, a model trained with data collected under one condition setting may not be able to generalize to other settings. This problem deters models from being deployed in practical scenarios. Aimed at making models more generalized and overcoming the domain shift problem, transfer-learning-based algorithms regarding either domain adaptation/generalization or zero/few-shot learning are also proposed [1], [31], [32], [33], [34], [35]. However, it is still an open problem regarding how to design a sophisticated transfer learning algorithm for FAS by considering all possible capturing settings.

III. METHODOLOGY

In this work, we propose a two-branch framework based on CNN, RNN inspired by how human beings can act to observe and explore spoofing clues. The overview of the framework

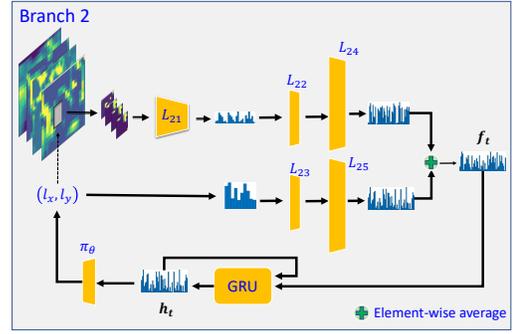


Fig. 4. Illustration of Branch 2. It consists of a Gated Recurrent Unit (GRU) and several linear layers (π_θ , L_{21} , L_{22} , etc.). It takes T steps for local features to be extracted. For each step, the reinforcement learning agent predicts a location (l_{x_t}, l_{y_t}) via its policy model (π_θ). Subsequently, at (l_{x_t}, l_{y_t}) , a sub-patch from F will be cropped as $F(l_{x_t}, l_{y_t})$, and $F(l_{x_t}, l_{y_t})$ will be further processed by the GRU. After T steps, the output hidden state of the GRU (h_T) serves as the extracted local features to be fused with the global features for the classification purpose.

is shown in Fig. 3. The backbone network firstly embeds the RGB image X into a feature map F . Then, F is forwarded to the subsequent Branch 1 and Branch 2 for the extraction of global and local discriminative features, respectively. Finally, these features are fused for the classification purpose in the final fully-connected (FC) layer. The details of Branch 2 are shown in Fig. 4, where L_{21} is a convolutional layer, and L_{22} , L_{23} , L_{24} and L_{25} are fully-connected layers. In Section III-A, we describe how the backbone and Branch 1 work cooperatively to extract global features. Subsequently, in Section III-B, we illustrate the extraction of local features from suspicious sub-patches sequentially. Afterward, the reinforcement learning leveraged to predict the locations of those sub-patches will be elaborated in Section III-C. Finally, in Section III-D, we present the optimization process.

A. Global Feature Extraction

The global feature extraction aims to exploit global discriminative information (e.g., paper boundaries, bezels, salient reflection patterns). In our framework, global features are extracted in Branch 1, a sequence of convolutional blocks. Branch 1 processes F from the backbone to extract global features f_g , into which all elements from F are encoded.

For implementation, we construct the backbone and Branch 1 based on the ResNet18 [36] architecture such that we can fairly compare our method with the recent ResNet-based methods (e.g., [24]). In particular, we adopt the first convolutional layer and four subsequent residual blocks of ResNet18 as our backbone network. The remaining convolutional residual blocks and the Global Average Pooling (GAP) layer [37] constitute Branch 1. The details of these modules are provided in Table I.

B. Local Feature Extraction

The local features are expected to exploit discriminative information from a local patch. The local feature extraction consists of T steps in total. We first elaborate on the procedure of local feature extraction at a certain step t . At the beginning of step t , a location (l_{x_t}, l_{y_t}) is first produced, where

l_{x_t} and l_{y_t} represent the horizontal and vertical coordinates respectively. The learning of the location prediction is via reinforcement learning, which will be introduced in the next sub-section. Then, from F , a square sub-patch $F^{(l_{x_t}, l_{y_t})}$ with the patch size p centering at (l_{x_t}, l_{y_t}) is cropped. Next, $F^{(l_{x_t}, l_{y_t})}$ together with the location information (l_{x_t}, l_{y_t}) will be encoded to an intermediate feature f_t . As such, f_t contains information from the observation at step t .

While previous works usually utilize RNN to leverage temporal information from sequential video frames, we particularly employ a Gated Recurrent Unit (GRU [27]) to learn local features from f_1, f_2, \dots, f_t in a sequential and recursive manner. The reason is that the hidden state h_t of the GRU can be learned from h_{t-1} and f_{t-1} :

$$\begin{aligned} z_t &= \text{sigmoid}(W_z f_t + U_z h_{t-1} + b_z) \\ q_t &= \text{sigmoid}(W_q f_t + U_q h_{t-1} + b_q) \\ \hat{h}_t &= \tanh(W_h f_t + U_h (q_t \odot h_{t-1}) + b_h) \\ h_t &= z_t \odot h_t + (1 - z_t) \odot \hat{h}_t, \end{aligned} \quad (1)$$

where \odot is the Hadamard product operation, $W_{\{z,q,h\}}$, $U_{\{z,q,h\}}$, and $b_{\{z,q,h\}}$ are parameters of the GRU. Furthermore, when analyzing a presented example, human beings' knowledge with respect to the example grows as information is gradually gained at each step. In other words, based on previous observations (f_1, f_2, \dots, f_t) , their assessment toward a suspicious example can get reinforced after a new observation (f_{t+1}) . Therefore, we specifically employ the GRU to memorize the observed information and learn local features. After T steps, h_T is treated as the final extracted local feature because it has perceived the local information during the T steps of observations. Finally, the proposed framework jointly exploits global and local features by fusing h_T and f_g for the classification purpose.

C. Reinforcement Learning for Face Anti-Spoofing

To explore spoofing-discriminative local information, we leverage reinforcement learning to train an agent that can help predict locations of sub-patches where spoofing clues may appear. In this context, a reinforcement learning agent is an abstract subject that explores clues in a certain environment and predicts locations.

Environment: In our framework, we treat F from the backbone as the environment where our agent predicts locations and gets feedbacks to update its policy. This is because the backbone can distill shallow spoofing-related features from raw RGB pixels into F . Thus, F can especially provide spoofing-related information for the agent to predict appropriate locations to extract effective local features.

Although an input RGB image (X) can also be set as the environment, we experimentally find that using the backbone to extract F and setting it as the environment can provide better results. We conjecture that raw pixels of an RGB image may contain interference such that the agent could get overwhelmed in a complex environment. On the other hand, the backbone could filter out unrelated information and distill spoofing-related information from raw pixels, and thus provide

a specific environment F for the agent to explore spoofing clues with less disturbance. The experimental results in IV-C1 show the superiority of using the backbone and setting F as the environment.

State: At step t , according to a certain policy, the agent predicts the location of a sub-patch based on its state s_t . As h_t carries history action information, we define the agent's state as h_t : $s_t = h_t$.

Action and Policy: In the framework, the agent learns to predict the location of a sub-patch to explore spoofing information. Hence, our agent's action a_t is to predict the location: $a_t = (l_{x_t}, l_{y_t})$. Then, the sub-patch at (l_{x_t}, l_{y_t}) will be cropped for the extraction of local features.

For effective location prediction, an optimal policy π should be provided to guide our agent to predict a location according to its current state: $a_t = \pi(s_t)$. Following the policy gradient theory [38], we parameterize π as π_θ by using a differentiable linear layer, where θ denotes its parameters. In this way, we can optimize θ with the standard backward propagation based on reward signals, which will be illustrated later (Section III-D).

Reward: After predicting the locations, the agent should get reward signals to evaluate how discriminative the information that the sub-patches contain for the classification. The more effective the predictions, the higher the rewards. Since the classification is conducted at the final step, we define a delayed reward as

$$r_t = \begin{cases} 0, & \text{if } t < T \\ \log P(y_{gt} | \mathbf{X}), & \text{if } t = T, \end{cases} \quad (2)$$

where r_t is the reward signal at step t , T is the total number of observation steps, y_{gt} is the ground-truth label of \mathbf{X} , and P is the predicted probability distribution over the binary labels. The agent will be trained to gain the cumulative reward

$$R = \sum_{t=1}^T r_t = \log P(y_{gt} | \mathbf{X}) \quad (3)$$

as high as possible.

D. Training and Optimization

1) *Two-stage training scheme:* When optimizing our framework, although end-to-end training is achievable, it may not provide satisfactory performance. As mentioned, the output feature map F from the backbone can be seen as an environment where our agent acts to learn its policy. If the backbone is involved in training, the environment will be unstable. Assume that the training is in epoch i , π_θ is optimized according to the "environment" F_i . However, if the backbone is also involved in optimization, the "environment" will change in the next epoch, meaning that $F_{i+1} \neq F_i$. Therefore, the agent may not act properly to the new "environment" F_{i+1} with the policy learned from the F_i .

To tackle this problem, we propose to use a two-stage training scheme. At the first stage, we pretrain a ResNet18 model with the training data. Then, the parameters of the corresponding modules will be loaded to the backbone from the pretrained model. Subsequently, the parameters of the backbone will be frozen and not involved in the second-stage optimization such that $F_{i+1} = F_i$. As such, fixing the

parameters of the backbone is to fix F , which can help keep a stable “environment” and extract more effective local features. The experimental results in Section IV-C show the superiority of our two-stage training.

2) *Joint optimization*: At the second stage, we optimize the parameters of Branch 1 and 2 jointly. The parameters of Branch 1 and Branch 2 except π_θ are optimized by the standard cross-entropy loss with binary labels for supervision.

As for π_θ , it is optimized with reinforcement learning. The optimization of π_θ can be formulated as the maximizing of the following objective function:

$$J(\theta) = \mathbb{E}_{\rho(s_{1:T};\theta)}[R], \quad (4)$$

where $\rho(s_{1:t};\theta)$ is the distribution over action sequences, which depends on π_θ .

According to the policy gradient theory [38], we adopt a differentiable linear layer to approximate the policy function. Hence, the maximization of $J(\theta)$ can be via the calculation of the gradient of $J(\theta)$ and the application of the gradient ascend. To this end, we leverage the REINFORCE rule [39] to approximate the gradient of $J(\theta)$:

$$\begin{aligned} \nabla_\theta J(\theta) &= \sum_{t=1}^T \mathbb{E}_{\rho(s_{1:t};\theta)}[\nabla_\theta \log \pi_\theta(a_t|s_{1:t};\theta)R] \\ &\approx \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t|s_{1:t};\theta)R. \end{aligned} \quad (5)$$

As the gradient of $J(\theta)$ can be simply computed by $\nabla_\theta \log \pi_\theta(a_t|s_{1:t};\theta)$. Thus, π_θ can be optimized by the standard backward propagation with this approximated gradient.

IV. EXPERIMENTS

This section describes how we conduct experiments to evaluate our method. To begin with, we introduce six benchmark databases employed in the experiments. After that, we illustrate the implementation details. Finally, we present and discuss the experimental results.

A. Databases

We utilize six publicly available face presentation attack databases in our experiments, including CASIA Face Anti-Spoofing Database [40], IDIAP REPLAY-ATTACK [3], MSU Mobile Face Spoofing Database [17], OULU-NPU database [41], the Spoofing in the Wild (SiW) database [9] and the ROSE-YOUTU database [1].

1) *CASIA FASD*: The CASIA Face Anti-Spoofing Database (CASIA for short) has 20 and 30 subjects in its training and testing set respectively. For each subject, there are 12 videos, among which 3 genuine face videos are recorded from the genuine faces and 9 spoofing face videos from photos and screens. As such, the CASIA database has 600 videos in total.

2) *IDIAP REPLAY-ATTACK*: The IDIAP REPLAY-ATTACK database [3] (IDIAP for short) is constituted of 1,200 videos in total, with 360, 360 and 480 videos in the training set, development set, and testing set, respectively. In this database, there are two illumination conditions: 1) a

controlled condition where the background is uniform and the source of lighting is a fluorescent lamp; 2) an adverse condition where the background is non-uniform and the source of lighting is daylight. The attack videos of each subject involve the 1) Print Attack: High-resolution face pictures printed on a paper. 2) Replay Attack: High-resolution pictures or videos were displayed on the screen of an iPhone 3GS and an iPad. To collect such data, the webcam of a MacBook, an iPhone 3GS and a Canon PowerShot camera are used.

3) *MSU MFSD*: The MSU Mobile Face Spoofing Database (MSU for short) [17] consists of 280 video clips of photo and video attack from 35 subjects. Two types of cameras are used to collect videos: the built-in camera in MacBook Air 13 (640 × 480) and the front-facing camera in the Google Nexus 5 Android phone (720 × 480). However, all the videos are only collected in normal indoor lighting environments.

4) *OULU-NPU*: The OULU-NPU database, similar to the IDIAP database, is divided into the training, development, and testing set with 20, 15, and 20 subjects, respectively. Overall, it contains 4950 face videos collected under three different environment conditions (e.g., different illumination and background conditions), with the frontal cameras of six mobile phones. As for attack mediums, two printers and two display devices were used to produce print attack and video attack examples. Furthermore, the OULU-NPU database provides four protocols for evaluation. Among them, Protocol 1, 2, and 3 aim to evaluate a model’s generalization capability to unseen environment conditions, unseen attack mediums, and unseen camera modules, respectively. Protocol 4 simultaneously considers the unseen environment conditions, attack mediums, and camera modules.

5) *SiW*: The Spoofing in the Wild (SiW for short) database [9] covers 165 subjects. For each subject, eight genuine face videos and 20 spoofing face videos are recorded. As for data collection environments, four sessions with variations of distances, poses, illuminations, and expressions have been considered [9]. For print attack examples, an HP Color LaserJet M652 printer is for printing high resolution (5184 × 3456) and low resolution (1920 × 1080) photos. Besides, four e-devices (Samsung Galaxy S8, iPhone 7, iPad Pro, and PC ASUS MB168B) are used to collect spoofing faces on their screens. As for cameras, a Canon EOS T6 and a Logitech C920 webcam are utilized to capture data. Totally, the SiW database has up to 4478 genuine and spoofing face videos. Also, it offers three protocols to evaluate the generalization capability of a model to unseen face poses and expressions (Protocol 1), unseen attack producing mediums (Protocol 2), and unseen Presentation Attack types (Protocol 3).

6) *ROSE-YOUTU*: The ROSE-YOUTU Face Liveness Detection Database (ROSE-YOUTU for short) is collected by the industry partner, YouTu. It involves 20 subjects, and for each subject, there are 25 genuine and 150 to 200 spoofing face videos. The data is diversely collected, covering up to 5 different lighting conditions. Also, the front-facing cameras of five mobile devices (Hasee phone, Huawei phone, ZTE phone, iPad, iPhone 5s) were used to record the videos, with resolution ranging from 640 × 480 to 1280 × 720. Moreover, besides

TABLE I

THE PARAMETERS OF EACH MODULE IN THE FRAMEWORK. “2D CONV” DENOTES THE A SEQUENCE OF A 2D CONVOLUTIONAL LAYER, A 2D BATCH NORMALIZATION LAYER AND A ReLU LAYER. “LINEAR” DENOTES A SEQUENCE OF A FULLY-CONNECTED LAYER AND A ReLU LAYER. “GAP” DENOTES THE GLOBAL AVERAGE POOLING LAYER. THE INPUT AND OUTPUT SIZES OF EACH LAYER ARE SHOWN. p IS THE PATCH SIZE FOR THE CROPPING.

Module	Layer	Type	Size
Backbone	1	2D Conv	In: $3 \times 256 \times 256$ Out: $64 \times 256 \times 256$
	2	MaxPooling	In: $64 \times 256 \times 256$ Out: $64 \times 128 \times 128$
	3	Residual Block $\times 2$	In: $64 \times 128 \times 128$ Out: $64 \times 64 \times 64$
	4	Residual Block $\times 2$	In: $64 \times 64 \times 64$ Out: $128 \times 32 \times 32$
Branch 1	1	Residual Block $\times 2$	In: $128 \times 32 \times 32$ Out: $256 \times 32 \times 32$
	2	Residual Block $\times 2$	In: $256 \times 32 \times 32$ Out: $512 \times 32 \times 32$
	3	GAP	In: $512 \times 32 \times 32$ Out: $512 \times 1 \times 1$
Branch 2	L_{21}	2D Conv + GAP	In: $128 \times p \times p$ Out: $256 \times 1 \times 1$
	L_{22}	Linear	In: 256 Out: 512
	L_{23}	Linear	In: 2 Out: 512
	L_{24}	Linear	In: 512 Out: 512
	L_{25}	Linear	In: 512 Out: 512
	GRU	GRU	In: 512 Out: 512
	π_θ	Linear	In: 512 Out: 512
FC	1	Linear	In: 1024 Out: 2

(still and quivering) printed photos and replay video examples (displayed with Lenovo LCD screen and Mac screen), the ROSE-YOUTU database further includes various paper mask attack examples. Such attacks can contain 3D information but are lacked in the aforementioned five databases. Hence, we leverage the ROSE-YOUTU database to evaluate our method further.

B. Experiments Settings

1) *Evaluation protocols and metrics*: When evaluating the proposed framework, we report the experimental results in terms of Equal Error Rate (EER), Half-Total Error Rate (HTER), Attack Presentation Classification Error Rate (APCER), Bona Fide Presentation Classification Error Rate (BPCER), and Average Classification Error Rate (ACER) in different scenarios. For intra-database experiments on the CASIA, IDIAP, and ROSE-YOUTU databases, we use data of the training set of the given database to train models and report EER results on the corresponding testing sets. When conducting cross-database experiments, we report HTER. Besides, for the OULU-NPU and SiW databases, we respectively follow the four protocols [41] and the three protocols [9] to evaluate the generalization capability of our method by reporting ACER, APCER, and BPCER results.

2) *Implementation Details*: As for the framework input, we consider including background information as spoofing-related

TABLE II

PERFORMANCE COMPARISONS BETWEEN THE MODELS WITH ONLY GLOBAL FEATURES (BRANCH 1), WITH ONLY LOCAL FEATURES (BRANCH 2) AND WITH BOTH FUSED ON THE CASIA, ROSE-YOUTU, REPLAY-ATTACK AND THE FOUR PROTOCOLS (P1, P2, P3, P4) OF OULU-NPU DATABASES. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Methods	Only local features	Only global features	Fused (ours)
CASIA	3.26	0.66	0.17
IDIAP	2.54	0.88	0.00
ROSE-YOUTU	8.06	5.42	1.79
OULU-NPU-P1	6.70	4.06	2.58
OULU-NPU-P2	5.33	3.08	1.15
OULU-NPU-P3	3.38 ± 3.78	2.21 ± 1.48	1.18 ± 1.19
OULU-NPU-P4	11.23 ± 5.94	5.07 ± 1.78	3.12 ± 2.01

clues are diverse and may not necessarily appear on face areas, which can be implemented by expanding face detection bounding boxes to crop out faces. However, there are various attack scenarios, and optimal bounding box sizes could depend on scenarios. Since the entire video frames can be regarded as the detected face cropped by the bounding box of a special configuration, we use such configuration by default for the framework input as a consistent way to evaluate the proposed framework. Nevertheless, we also evaluate the effectiveness of the framework under different configurations of bounding boxes, where the framework inputs are consistently resized to 256×256 pixels. Then, the backbone network embeds the input into feature maps $\mathbf{F} \in \mathbb{R}^{128 \times 32 \times 32}$. Subsequently, \mathbf{F} is forwarded to Branch 1 and Branch 2 to extract global and local features respectively. For the final classification, we fuse the global and local features from Branch 1 and 2 by using the Concatenation as the input of the final Fully Connected (FC) layer. We show the details of the backbone network, Branch 1, Branch 2, and the final FC layer in Table I, where p denotes the size for cropping patches.

When training the framework, we follow the two-stage training scheme stated in Section III. At the first stage, we pretrain a ResNet18 model [36] with cross-entropy loss with training data. The trained parameters of the first convolutional layer and the four subsequent residual blocks are then loaded to the backbone, and the parameters of the backbone will be fixed and excluded from the second-stage optimization. At the second stage, the GRU’s hidden state \mathbf{h}_t is initialized as \mathbf{f}_g , and the location of the initial patch is sampled from a normal distribution whose symmetry center corresponds to the center of the input images. Then, Branch 1 and 2 will be optimized jointly from scratch with the standard backward propagation with gradients of the cross-entropy loss and $J(\theta)$. By default, except for the declaration, the input configuration is “FULL”; the number of observation steps T is set as 8; the patch size p is set as 8; and the fusion method is set as the Concatenation. In addition, we also explore the impacts of p , T , and different fusion methods in Section IV-C.

C. Experimental results

1) *Analysis of Jointly Using Global and Local Features*: In this subsection, we demonstrate the effectiveness of our

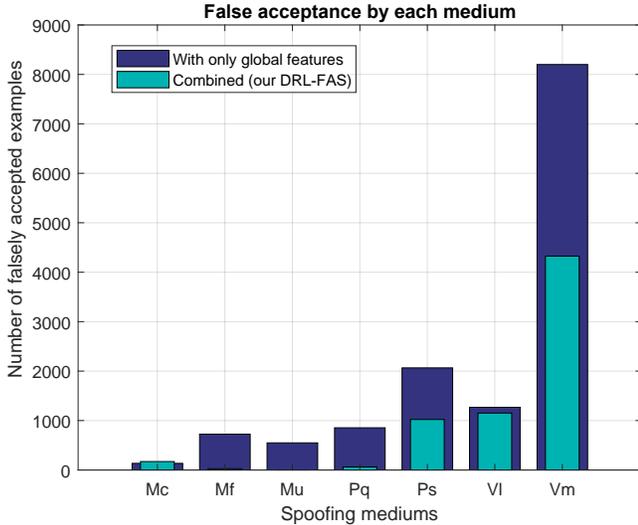


Fig. 5. The statistics of the false acceptance examples on the ROSE-YOUTU database [1]. The horizontal axis denotes the presentation attack mediums, and the vertical axis indicates the number of falsely accepted examples. “Mc” indicates a paper mask with two eyes and mouth cropped out. “Mf” indicates a paper mask without cropping. “Mu” indicates a paper mask with the upper part cut in the middle. “Pq” indicates a quivering printed paper. “Ps” indicates a still printed paper. “VI” indicates a video that records a Lenovo LCD display. “Vm” indicates a video that records a Mac LCD display. Some examples of “Vm” and “VI” can be seen in Fig. 1.

proposed framework by jointly utilizing global and local features. To this end, we ablate Branch 1 and Branch 2 separately to compare results with only local features and with only global features. As shown in Table II, the results with only global features (Branch 2 ablated) are better than the results with only local features (Branch 1 ablated), which indicates that global features can be more effective than local features. Intuitively, people are likely to provide a reliable assessment with a glance at the original video frame, especially when discriminative artifacts appear, e.g., paper boundaries, bezels, or obvious reflections. However, when given merely a few local sub-patches, human beings could have trouble assessing the liveness of the original examples as discriminative artifacts may be absent in these patches. This is just like what the story *Blind Men and An Elephant* [42] tells. Moreover, by fusing the global and local features, our proposed framework can further achieve better performance. Such improvement supports our motivation that “taking closer observations at local sub-patches” can provide more information to refine the classification.

For further analysis, we collect the statistical results of the falsely accepted examples for each medium on the ROSE-YOUTU database. Fig. 5 shows that, with only global features, the number of falsely-accepted “Vm” (video attack recorded from a Mac display) is nearly 9000, which represents the largest proportion. By carefully reviewing the data in the ROSE-YOUTU database, we find that the visual quality of “Vm” is generally better than that of “VI” (video attack recorded from a Lenovo display) as the resolution of a Mac display (2560×1600 resolution) is much higher than that of a Lenovo LCD (1920×1080 resolution). For instance, Fig. 1(c) and is a “VI” example, and Fig. 1(e) is a “Vm”



Fig. 6. Illustrations of the faces cropped by bounding boxes of different configurations. (a) is the entire original video frame of a spoofing example from the OULU-NPU database, which can be regarded as the face image cropped by a bounding box with a special configuration (“FULL”). (b) is the detected face cropped from (a) by a bounding box with the default setting of the dlib’s CNN detector ($\alpha = 0.0$). Analogously, (c) is that by 20% ($\alpha = 0.2$), (d) is that by 40% ($\alpha = 0.4$), and (e) is that by 60% ($\alpha = 0.6$).

TABLE III
THE COMPARISON BETWEEN THE RESNET18 BASELINE AND OUR PROPOSED FRAMEWORK. THE ACER (%) RESULTS ARE REPORTED ON THE FOUR PROTOCOLS (P1, P2, P3, P4) OF THE OULU-NPU DATABASE. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Method	ACER			
	P1	P2	P3	P4
ResNet18 ($\alpha = 0.0$)	4.0	4.5	7.2±7.0	10.4±6.0
ResNet18 ($\alpha = 0.2$)	5.3	3.3	4.8±5.0	9.4±7.1
ResNet18 ($\alpha = 0.4$)	2.4	3.6	4.5±4.1	7.4±3.8
ResNet18 ($\alpha = 0.6$)	2.3	3.7	5.1±2.4	6.3±2.6
ResNet18 (FULL)	7.2	2.3	3.4±1.3	11.4±5.1
Ours ($\alpha = 0.0$)	3.7	4.0	6.4±6.9	9.5±6.6
Ours ($\alpha = 0.2$)	4.4	3.0	3.9±2.3	8.6±5.6
Ours ($\alpha = 0.4$)	1.6	3.4	2.9 ±1.4	6.0±4.6
Ours ($\alpha = 0.6$)	1.4	2.6	3.1±1.8	5.0 ±3.7
Ours (FULL)	4.7	1.9	3.0±1.5	7.2±3.9

example. As shown, the visual quality of the “Vm” looks better and than the “VI”. In other words, the spoofing faces of “Vm” visually look more similar to genuine faces than “VI” (as well as others). Therefore, “Vm” examples get most falsely accepted as genuine faces than “VI” and the others. Moreover, as shown in Fig. 5, our method can generally lead to fewer falsely accepted spoofing examples of various attack mediums. Although the falsely accepted “Vm” from our method still accounts for the largest proportion of false acceptance, the value is nearly half of that with only global features. This means that our method can better discriminate spoofing examples of good visual quality by leveraging local information from sub-patches. It corresponds to our motivation that people can zoom in local sub-patches and explore subtle spoofing clues to refine and confirm their assessment of the liveness of examples. By far, we demonstrate that our framework can better counter spoofing attacks by jointly using global and local features. In the next subsection, we investigate how well our framework can be generalized to inputs of other configurations.

2) *Analysis of Configurations of the Framework Input*: As our framework is proposed to exploit the discriminative information which may not necessarily appear on face areas, we also propose to investigate the performance by configuring the input with different scales of information from backgrounds based on detected faces. To this end, we propose to use a

TABLE IV

THE COMPARISON OF DIFFERENT METHODS FOR SELECTING PATCHES. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

EER	MAX-SCORES	RANDOM	DRL (ours)
CASIA	0.25	0.32	0.17
IDIAP	0.00	0.00	0.00
ROSE-YOUTU	3.29	2.14	1.79
OULU-NPU-P1	3.82	3.57	2.58
OULU-NPU-P2	1.54	1.76	1.15
OULU-NPU-P3	1.20±1.03	1.38±0.94	1.18±1.12
OULU-NPU-P4	4.67 ±1.31	4.32±1.75	3.12±2.01

dlib's CNN face detector [43] to obtain detection bounding boxes. Subsequently, a bounding box will be expanded by 0%, 20%, 40%, and 60% (i.e., $\alpha = \{0.0, 0.2, 0.4, 0.6\}$) to produce four face images that have different scales of background information. Besides, using the entire video frames can be treated as cropping faces with a special configuration for the bounding box, and we denote such configuration as "FULL". As such, there will be five groups of face images in different configurations for the framework input, and some of the examples are shown in Fig. 6. For experiments, we adopt the OULU-NPU database, as it provides four protocols for extensive evaluation. Also, we train ResNet18 models to provide the baseline results, where only global features are considered. The experimental results are shown in Table III. It is obvious that our method can achieve better ACER results than the counterparts of ResNet18 over different input configurations. To sum up, the experiments show that the proposed method can still be effective in jointly using global and local features extracted from face images that have different scales of background information.

3) *Effect of using reinforcement learning*: In this subsection, we show the effectiveness of adopting deep reinforcement learning (DRL) for selecting patches. To be more specific, we compare our proposed DRL with the method of selecting patches that have the max SoftMax scores (denoted as the MAX-SCORES method), and the method of selecting patches randomly (denoted as the RANDOM method). In the implementation of the MAX-SCORES method, we pretrained a patch-based CNN based on [29] with the training data to infer the SoftMax scores of the patches, and those patches that have the max SoftMax scores are selected for the framework. Besides, in the implementation of the RANDOM method, we use a random number generator to generate locations of the patches to be selected. Table IV compares the performance of the MAX-SCORES, RANDOM, and the DRL with respect to patch selection. Regarding the IDIAP database, the MAX-SCORES, RANDOM, and DRL methods achieve 0.00% EER, meaning that local features can help with the final prediction, regardless of how we select patches. Nevertheless, in the other experiments based on the ROSE-YOUTU, CASIA, and OULU-NPU databases, our proposed DRL shows the effectiveness in selecting patches by achieving the best EER results.

TABLE V

THE EER RESULTS OF DIFFERENT PATCH SIZE p ON THE CASIA, ROSE-YOUTU AND PROTOCOL 1 OF THE OULU-NPU DATABASE (OULU-NPU-P1). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD, AND THE SECOND-BEST RESULTS ARE UNDERLINED. IN THE EXPERIMENTS, WE SET $T = 8$.

EER	CASIA	ROSE-YOUTU	OULU-NPU-P1
$p = 2$	0.35	3.92	2.29
$p = 4$	0.17	<u>2.73</u>	2.95
$p = 8$	0.17	1.79	<u>2.58</u>
$p = 16$	<u>0.28</u>	3.65	4.20

TABLE VI

THE EER RESULTS OF DIFFERENT TOTAL NUMBER OF OBSERVATION STEPS T ON THE CASIA, ROSE-YOUTU, AND PROTOCOL 1 OF THE OULU-NPU DATABASE (OULU-NPU-P1). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD, AND THE SECOND-BEST RESULTS ARE UNDERLINED. IN THE EXPERIMENTS, WE SET $p = 8$.

EER	CASIA	ROSE-YOUTU	OULU-NPU-P1
$T = 2$	0.18	1.64	3.85
$T = 4$	<u>0.17</u>	2.13	<u>3.15</u>
$T = 8$	<u>0.17</u>	<u>1.79</u>	2.58
$T = 16$	0.06	1.92	3.52

4) *Analysis of Local Features Extraction*: In this subsection, we analyze the impact of patch size p and the number of steps T for the local feature extraction.

Effect of patch size p To analyze the effect of p , we conduct experiments with $p = \{2, 4, 8, 16\}$, and we provide the EER results on the CASIA, ROSE-YOUTU, and the OULU-NPU-P1 in Table V. Regarding the CASIA database, we can see that the EER results become better from 0.35% to 0.17% when p increases from 2 to 8 but deteriorate when we further increase p to 16 (0.28%). Regarding the ROSE-YOUTU database, a similar trend can be observed that the EER performance improves up to 1.79% EER when p increases from 2 to 8. However, when $p = 16$, the EER performance drops as 3.65% EER. Regarding the OULU-NPU-P1, the best EER is achieved when $p = 2$, and the EER result of $p = 8$ is better than that of $p = 4$ and $p = 16$. Therefore, the optimal p is different for different databases, and simply increasing p may not necessarily lead to better performance. Nevertheless, we observe that $p = 8$ can achieve the desired performance in general, and thus we fix $p = 8$ for all other experiments.

Effect of total number of observation steps T Besides the size of the local patches, we are also interested in the impact of different numbers of observation steps T . To this end, we conduct experiments by increasing T from 2 to 16, and the results are reported in Table VI. As we can observe, for the CASIA database, the EER performance improves when T increases from 2 to 16. Regarding the ROSE-YOUTU database, the best EER result of 1.65% is achieved when $T = 2$, and the second-best EER of 1.79% is achieved when $T = 8$. Regarding the OULU-NPU-P1, the EER performance improves as T changes from 2 to 8 but gets worse when $T = 16$. In summary, simply increasing T does not necessarily lead to better performance. We conjecture the reason that when T is increased to include more patches, those patches

TABLE VII

THE COMPARISON BETWEEN WITH/WITHOUT THE BACKBONE FOR THE FEATURE EMBEDDING. THE EXPERIMENTS ARE CONDUCTED ON THE CASIA, ROSE-YOUTU, REPLAY-ATTACK, OULU-NPU DATABASES. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Methods	Without backbone	With backbone
CASIA	7.80	0.17
IDIAP	1.13	0.00
ROSE-YOUTU	2.38	1.79
OULU-NPU-P1	8.43	2.58
OULU-NPU-P2	1.04	1.15
OULU-NPU-P3	1.15 \pm 1.26	1.18 \pm 1.19
OULU-NPU-P4	9.57 \pm 5.35	3.12 \pm 2.01

TABLE VIII

PERFORMANCE COMPARISONS BETWEEN THE ONE-STAGE END-TO-END TRAINING AND OUR TWO-STAGE TRAINING ON THE CASIA DATABASE. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%). IN THE EXPERIMENTS, WE SET $p = 8$.

Observation steps	One-stage	Our two-stage
$T = 2$	20.1	0.184
$T = 4$	18.7	0.171
$T = 8$	4.32	0.171

containing less discriminative information may deteriorate the performance. Nevertheless, we consistently choose $T = 8$ for other experiments as it can lead to the best or second-best performance in Table VI.

5) *Analysis of other settings of the framework:* In this subsection, we analyze the impacts of the backbone, training strategies, and feature fusion methods in our proposed framework.

Effect of using the backbone for local features To study the effectiveness of using the backbone for feature embedding and local feature extraction, we also implement the framework without the backbone for feature embedding as a baseline, where patches are cropped at the predicted locations from original RGB inputs instead of the embedded feature maps. Table VII shows the results for comparison. We observe that better performances can be achieved in general by considering a backbone network. Such improvement indicates that extracting local features from the feature maps through the backbone can help to extract more discriminative spoofing-related information. Moreover, the results with the backbone on the OULU-NPU-P1 and -P4 (2.58% and 3.12%) are significantly better than those without the backbone (8.43% and 9.57%). As Protocol 1 and 4 involve unseen environments in the testing data, conducting feature embedding through the backbone network to extract local features may alleviate environmental interference to some extent such that better performance can be achieved.

Effect of the two-stage training Although training our framework in an end-to-end manner is achievable, we propose to use a two-stage training scheme to better optimize our framework. Table VIII compares the results of the one-stage

TABLE IX

THE RESULTS ON THE CASIA DATABASE IN DIFFERENT SETTINGS OF FUSION METHODS. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%).

Observation steps	Patch size	Average	Weighted Average	Concatenation
$T = 4$	$p = 4$	16.1	12.3	0.184
	$p = 8$	9.54	11.6	0.171
$T = 8$	$p = 4$	15.6	12.4	0.171
	$p = 8$	0.132	0.172	0.171

TABLE X

INTRA-DATABASE EXPERIMENTS ON THE CASIA DATABASE AND THE IDIAP DATABASE. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%) AND HTER (%). “-” MEANS THE RESULT IS NOT AVAILABLE. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Method	IDIAP		CASIA
	EER	HTER	EER
CNN [7]	6.1	2.1	7.4
Color-LBP [4]	0.4	2.9	6.2
Bottleneck feature fusion + NN [23]	0.8	0.0	5.8
MSR-ResNet [24]	0.2	0.4	3.1
DRL-FAS (ours)	0.0	0.0	0.2

end-to-end training and our proposed two-stage training. The EER results of our two-stage training are all remarkably lower than 0.2% for $T = 2, 4, 8$. However, for the one-stage training experiments, when $T = 2$, the EER result is up to 20.1%. Although the EER decreases as T increases, the best EER result is still above 4% ($T = 8$), which is much higher than all the results achieved by our two-stage training. Therefore, Table VIII shows that our two-stage training can help achieve better results by providing a stable environment such that the agent can learn to extract effective features, even when $T = 2$.

Effect of fusion methods In this framework, the global and local features are fused for classification. In Table IX, we compare results among three different fusion methods, the Average, Weighted Average [24], and Concatenation. We observe that the EER results are all lower than 0.2% when $T = 8$ and $p = 8$. However, when $T < 8$ or $p < 8$, the EER results of the Average and the Weighted Average are higher than 9%, while the result with only global features is less than 1% (shown in Table II). We conjecture the reason that the Average and Weighted Average fuse global and local features by averaging the elements at each dimension. When $T < 8$ or $p < 8$, the extracted local features may not be effective enough. As a result, discriminative information contained by global features may be distorted after the average operation. However, when the extracted local feature is not representative enough, the Concatenation could maintain original global information to a larger extent. Therefore, the Concatenation can provide stable results (all lower than 0.2% EER) under different settings of T and p , and we fix the Concatenation as the fusion method in other experiments.

6) *Intra-database experiments:* For further evaluation, we conduct experiments on six benchmark databases and compare our proposed method with state-of-the-art methods.

TABLE XI

THE INTRA-DATABASE EXPERIMENT RESULTS ON THE ROSE-YOUTU DATABASE COMPARED WITH THE STATE-OF-THE-ART METHODS. THE PERFORMANCE IS EVALUATED IN TERMS OF EER (%). THE BEST RESULT IS HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Method	ROSE-YOUTU
	EER
CoALBP (YCBCR) [2]	17.1
CoALBP (HSV) [2]	16.4
AlexNet [1]	8.0
3D-CNN [31]	7.0
DeSpoofing [30]	12.3
DRL-FAS (ours)	1.8

TABLE XII

COMPARISON BETWEEN THE PROPOSED FRAMEWORK AND STATE-OF-THE-ART METHODS ON THE SiW DATABASE. THE PERFORMANCE IS EVALUATED IN TERMS OF ACER (%). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $T = 8$ AND $p = 8$.

Protocol	Method	ACER
1	Auxiliary [9]	1.00
	STASN [11]	1.00
	STASN+ [11]	0.30
	DRL-FAS (ours)	0.00
2	Auxiliary [9]	0.57±0.69
	STASN [11]	0.28±0.05
	STASN+ [11]	0.15±0.05
	DRL-FAS (ours)	0.00±0.00
3	Auxiliary [9]	8.31±3.81
	STASN [11]	12.10±1.50
	STASN+ [11]	5.85±0.85
	DRL-FAS (ours)	4.51±0.00

- Protocol 1, 2, and 3 are for evaluating models' generalization capability to unseen face poses and expressions, unseen attack mediums, and unseen Presentation Attack types, respectively.
- For experiments of Protocol 1, the testing is only done once so there are no terms of standard deviation

Results on the CASIA database and IDIAP database

Table X provides the results of intra-database experiments on the CASIA and IDIAP databases. On the CASIA database, our method attains 0.17% EER, which is the best. The best performance on the IDIAP database can also be seen from the 0% EER on the development (DEV) set and the 0% HTER on the testing (TEST) set. On both the two benchmark databases, our method achieves the best performance and shows its effectiveness.

Results on the ROSE-YOUTU database Table XI compares our method with the baseline methods on the ROSE-YOUTU database. Our method can achieve the lowest EER (1.8%), while the second-best method 3D-CNN [31] merely achieves 7.0% EER. This shows our method's superiority. In addition, to evaluate how paper mask attacks in the ROSE-YOUTU database can fail depth-based methods, we implement the DeSpoofing method¹ [30] because it leverages depth information for training. However, it merely achieves 12.3% EER, which indicates that when encountering paper mask

attack examples that are with depth information, depth-based methods could lose efficacy. By contrast, our method can still perform favorably against the paper mask attack examples in the ROSE-YOUTU database.

Results on the SiW database Table XII shows the ACER results of our proposed framework and state-of-the-art methods on the SiW database. The Auxiliary method [9], with extra depth map and rPPG signals, attains 1.0%, 0.57%, and 8.31% ACER for Protocol 1, 2, and 3, respectively. Besides, the "STASN+" method [11] collects extra data outside the database for data augmentation and achieves 0.3%, 0.15% and 5.58% correspondingly. However, without extra data and auxiliary signals but only binary labels for supervision, our method can achieve the best results for the three protocols (0.00%, 0.00%, and 4.51%, respectively). In addition, for experiments of Protocol 2 and 3, our method can manage to get the smallest standard deviation, showing better stability. Moreover, the ACER results of all the listed methods for Protocol 3 are much higher than the results for Protocol 1 and 2. This indicates the setting of unseen presentation attack types is more challenging than unseen faces poses and expressions as well as unseen attack mediums.

Results on the OULU-NPU database Table XIII compares our proposed method with state-of-the-art ones. In Protocol 1, our method achieves 4.7% ACER, better than the "MSR-ResNet" method (5.9%) [24], which is also based on ResNet18 [36]. In Protocol 2, our method achieves the best ACER (1.9%). In Protocol 3, the Auxiliary method achieves the lowest 2.9% ACER, while our method achieves a very close ACER of 3.0%. In Protocol 4, our method achieves the second-best ACER of 7.2%. Overall, in terms of ACER, the DeSpoofing is better than the Auxiliary in Protocol 1 and 4, while the Auxiliary is better than the DeSpoofing in Protocol 2 and 3. This comparison shows that there may not be a method that is always optimal for all scenarios. Nevertheless, our method shows its effectiveness by achieving the best or the second-best ACER in Protocol 2, 3, and 4. Furthermore, according to Table III, using a proper setting of face cropping for the framework input can lead to better performance. Also, one can always further improve the framework with advanced neural networks and auxiliary information.

7) *Cross-database experiments:* We also conduct cross-database experiments to evaluate the generalization capability of our method to different data domains. For conciseness, "A → B" denotes an experiment where we run the training with the database "A" and run testing with the database "B".

Table XIV provides the cross-database experimental results among the CASIA, IDIAP, and MSU databases. In the experiments of CASIA → IDIAP and IDIAP → CASIA, the Auxiliary method [9] achieves the best results (27.6% and 28.4% HTER respectively). On the other hand, our framework achieves second-best HTER (28.4% and 33.2%). Among methods without auxiliary information, such as [24], our performance is the best. Also, in both the experiments of IDIAP → MSU and MSU → IDIAP, we implement the DeSpoofing method [30], and it outperforms the other baseline methods by achieving 33.2% and 27.8% HTER respectively. However, our method can significantly surpass it with much lower HTER results

¹<https://github.com/yaojieliu/ECCV2018-FaceDeSpoofing>

TABLE XIII

EXPERIMENT RESULTS FOR THE FOUR PROTOCOLS ON THE OULU-NPU DATABASE. THE PERFORMANCE IS EVALUATED IN TERMS OF APCER (%), BPCER (%) AND ACER (%) ON THE TESTING (TEST) SET. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

Protocol	Method	Test			Protocol	Method	Test		
		APCER	BPCER	ACER			APCER	BPCER	ACER
1	GRADANT [44]	1.3	12.5	6.9	2	GRADANT[44]	3.1	1.9	2.5
	DeSpoofing[30]	1.2	1.7	1.5		DeSpoofing[30]	4.2	4.4	4.3
	Auxiliary [9]	1.6	1.6	1.6		Auxiliary [9]	2.7	2.7	2.7
	MSR-ResNet [24]	5.1	6.7	5.9		MSR-ResNet [24]	7.6	2.2	4.9
	DRL-FAS (ours)	5.4	4.0	4.7		DRL-FAS (ours)	3.7	0.1	1.9
3	GRADANT [44]	2.6 ±3.9	5.0±5.3	3.8±2.4	4	GRADANT [44]	5.0 ±4.5	15.0±7.1	10.0±5.0
	DeSpoofing [30]	4.0±1.8	3.8±1.2	3.6±1.6		DeSpoofing [30]	5.1±6.3	6.1 ±5.1	5.6 ±5.7
	Auxiliary [9]	2.7±1.3	3.1±1.7	2.9 ±1.5		Auxiliary [9]	9.3±5.6	10.4±6.0	9.5±6.0
	MSR-ResNet [24]	3.9±2.8	7.3±1.1	5.6±1.6		MSR-ResNet [24]	11.3±3.9	9.7±4.8	9.8±4.2
	DRL-FAS (ours)	4.6±3.6	1.3 ±1.8	3.0±1.5		DRL-FAS (ours, FULL)	8.1±2.7	6.9±5.8	7.2±3.9

- Protocol 1, 2, and 3 are for evaluating a model’s generalization capability to unseen environment conditions, unseen attack mediums, and unseen camera modules, respectively. Protocol 4 extends the evaluation to unseen sessions, attack mediums and camera modules at the same time.
- For experiments of Protocol 1 and 2, the testing is only done once, so there are no terms of standard deviation.

TABLE XIV

INTER-DATABASE RESULTS BETWEEN THE CASIA, IDIAP AND MSU DATABASES. THE PERFORMANCE IS EVALUATED IN TERMS OF HTER (%). “–” MEANS THE RESULT IS NOT AVAILABLE. ON THE LEFT OF “→” IS THE DATABASE USED FOR TRAINING AND ON THE RIGHT FOR TESTING. WE MARK THE BEST RESULTS IN BOLD, AND THE SECOND-BEST RESULTS ARE UNDERLINED. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Method	CASIA→IDIAP	IDIAP→CASIA	MSU→IDIAP	IDIAP→MSU
LBP [45]	47.0	39.6	45.5	45.8
LBP-TOP [45]	49.7	60.6	46.5	47.5
Motion [45]	50.2	47.9	–	–
CNN [7]	48.5	45.5	37.1	48.6
Color LBP [4]	37.9	35.4	44.8	33.0
Color Texture [4]	30.3	37.7	33.9	34.1
Auxiliary [9]	27.6	28.4	–	–
DeSpoofing [30]	28.5	41.1	<u>33.2</u>	<u>27.8</u>
MSR-MobileNet [24]	30.0	33.4	–	–
MSR-ResNet [24]	36.2	34.7	–	–
DRL-FAS (ours)	<u>28.4</u>	<u>33.2</u>	29.7	15.6

TABLE XV

THE INTER-DATABASE EXPERIMENTS WHERE THE MODELS ARE TRAINED WITH DATA OF ROSE-YOUTU DATABASE AND TESTED ON THE CASIA AND IDIAP DATABASES. THE PERFORMANCE IS EVALUATED IN TERMS OF HTER (%). “*” MEANS THE RESULTS ARE WITH THE OUTLIER REMOVAL PROPOSED IN [1]. ON THE LEFT OF → IS THE DATABASE USED FOR TRAINING AND THE RIGHT FOR TESTING. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. IN THE EXPERIMENTS, WE SET $p = 8$ AND $T = 8$.

Method	ROSE-YOUTU→CASIA	ROSE-YOUTU→IDIAP
AlexNET without DA [1]	32.6	43.6
AlexNET with KMM [1]	31.6	43.6
AlexNET with SA [1]	35.0	38.5
AlexNET with KSA [1]	33.9	42.0
AlexNET with SA* [1]	30.7	36.2
AlexNET with KSA* [1]	30.1	38.8
DeSpoofing [30]	37.2	38.5
DRL-FAS (ours)	8.1	20.0

(29.7% and 15.6%), which is the best.

Table XV provides the experimental results of ROSE-YOUTU→CASIA and ROSE-YOUTU→IDIAP. In the experiment of ROSE-YOUTU→CASIA, our method can achieve the best 8.1% HTER. Both the ROSE-YOUTU and CASIA database include spoofing attacks that have bezels and paper boundaries observed. Hence, the proposed framework can capture such discriminative artifacts to achieve good performance. Besides, as for attack samples in the IDIAP database, there are few paper boundaries and bezels observed in the samples. However, in the experiment of ROSE-YOUTU→IDIAP, the proposed framework is still effective and achieves the best HTER (20.0%), at least 16% HTER significantly lower than the others. In summary, in the cross-database experiments, our proposed framework still shows effectiveness when the spoofing artifacts and backgrounds are from different data domains.

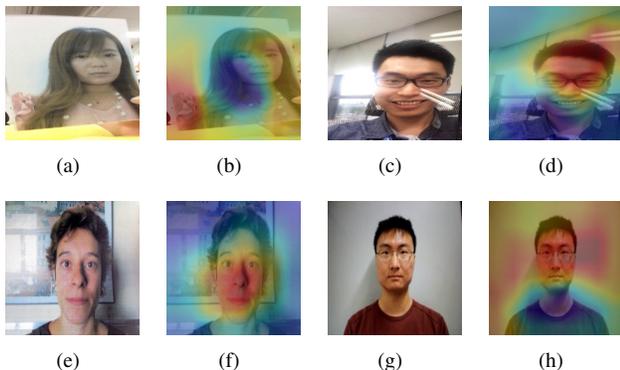


Fig. 7. The CAM heatmaps generated with global features. (a) is a print paper attack from the ROSE-YOUTU database, and its paper boundary is exposed on the left. (c) is a replay video attack from the ROSE-YOUTU database, and it has reflection patterns on the right. (e) and (g) are replay attack examples from the IDIAP REPLAY-ATTACK database. There are no discriminative artifacts. (b), (d), (f), and (h) are the CAM heatmaps of (a), (c), (e), and (g) respectively, where red/blue regions mean high/low activation (best viewed in color).

D. Visualization

In this subsection, we conduct further analysis based on visualization. To show what types of information the global features are likely to capture for anti-spoofing, we propose to apply the Class Activation Mapping (CAM) [46] to visualize activation heatmaps of global features. The CAM heatmaps are shown in Fig. 7, where red/blue areas mean high/low activation. Fig. 7(a) is a paper attack where the paper boundaries can be seen on the left and the right. Its CAM heatmap, Fig. 7(b), shows that the boundaries of both sides give high activation (red). Fig. 7(c) is a replay video attack with reflections appearing on the right, and we see from Fig. 7(d) that the reflection areas on the top right are red. Besides, we also explore the situation when the above artifacts are absent. Fig. 7(e) and Fig. 7(g) are replay video attack examples from the IDIAP database, and there are no discriminative artifacts observed. While the face area in Fig. 7(f) gives high activation, the background areas in Fig. 7(h) are also of high activation. In summary, discriminative information captured by global features may not necessarily appear on face areas, and information from backgrounds can also significantly contribute to anti-spoofing, even when bezels, reflections, etc., are not observed.

Besides, we propose to investigate how the local features can help with the performance by visualizing the predicted locations. To this end, we propose to fuse the global feature with the local feature h_t extracted at the step-index t ($t = 1, 2, \dots, 8$) for classification to get the confidence score c_t . In this way, the visualization results are shown in Fig. 8, where the number under each image is the confidence score c_t . As shown in Fig. 8, the performance of confidence scores shows an increasing trend. The first row of Fig. 8 is an attack example from the IDIAP database, a printed paper replayed in a video. As t increases, the predicted patches cover the printed stripes in the background, and c_t is generally increasing. The two rows below are printed paper attack examples from the CASIA database, where the boundaries of the printed paper can be treated as discriminative spoofing artifacts, and the

patches also cover the paper boundaries. In these two rows, the performance of c_t shows an increasing trend. The fourth row and fifth row are replay video attack examples from the OULU-NPU database and the ROSE-YOUTU database, respectively. We observe that the local patches mainly explore moir patterns. Meanwhile, there are $c_6 > c_7$ for the fourth and $c_3 > c_4$ for the fifth respectively, indicating that simply increasing t may not necessarily further improve the performance. Nevertheless, for $t > 1$, that $c_t > c_1$ still holds, showing that the information from patches can generally improve the overall performance.

Furthermore, we observe that the patches generally move from the center areas toward the boundary (background) areas. As the initial locations are sampled from a normal distribution whose symmetry center corresponds to the center of the input image, the initial patch is generally near the center areas. Thereby, the RNN extracts features from the center areas first. As spoofing features can also be found in the background areas, driven by reinforcement learning, the patches then move toward the background areas such that the RNN can extract features from these areas to improve performance.

We also visualize misclassified spoofing examples from the CASIA, ROSE-YOUTU, and OULU-NPU databases, which are shown in Fig. 9. Fig. 9(a) is a printed paper attack from the CASIA database. Although the paper boundary at the bottom can be seen, it is not obvious, and most of the other areas have no blur or distortion observed. Fig. 9(b) is a replay video attack example from the ROSE-YOUTU database, and Fig. 9(c) is a replay video attack example from the OULU-NPU database. These figures show few discriminative artifacts, such as reflection. Based on the observations from these examples, as there are few discriminative artifacts observed, the extracted global and local features may not be effective in differentiating the spoofing faces from genuine ones.

E. Computational Analysis

In this subsection, we analyze the computational efficiency of our proposed method and the ResNet18. As we can see from Table XVI, the total number of parameters of our model is about 16.50M, while the ResNet18 is 11.18M. This increased amount of parameters is reasonable as we introduce a local branch in our framework and our model can achieve better performance in the task of FAS. Despite that our model size increases by about 50% compared with the ResNet18, our proposed method does not introduce too much computational burden. In terms of Multiply-Accumulate-Operations (MACs) [47], which is for measuring the total multiplication and addition operations required for calculation, our method merely has more 0.04 Giga (1.7%) and 0.07 (3.0%) Giga when $T = 4$ and $T = 8$ than the ResNet18 baseline. Last but not least, we consider the inference efficiency by reporting Frames per Second (FPS) based on PyTorch and NVIDIA GTX 1080 Ti GPU. As we can see, while the ResNet18 achieves 150 FPS, our method can achieve 110 FPS and 70 FPS when $T = 4$ and $T = 8$, which is reasonable as the local branch works recurrently. Nevertheless, in practice, one can always use various neural network acceleration techniques to speed up models.

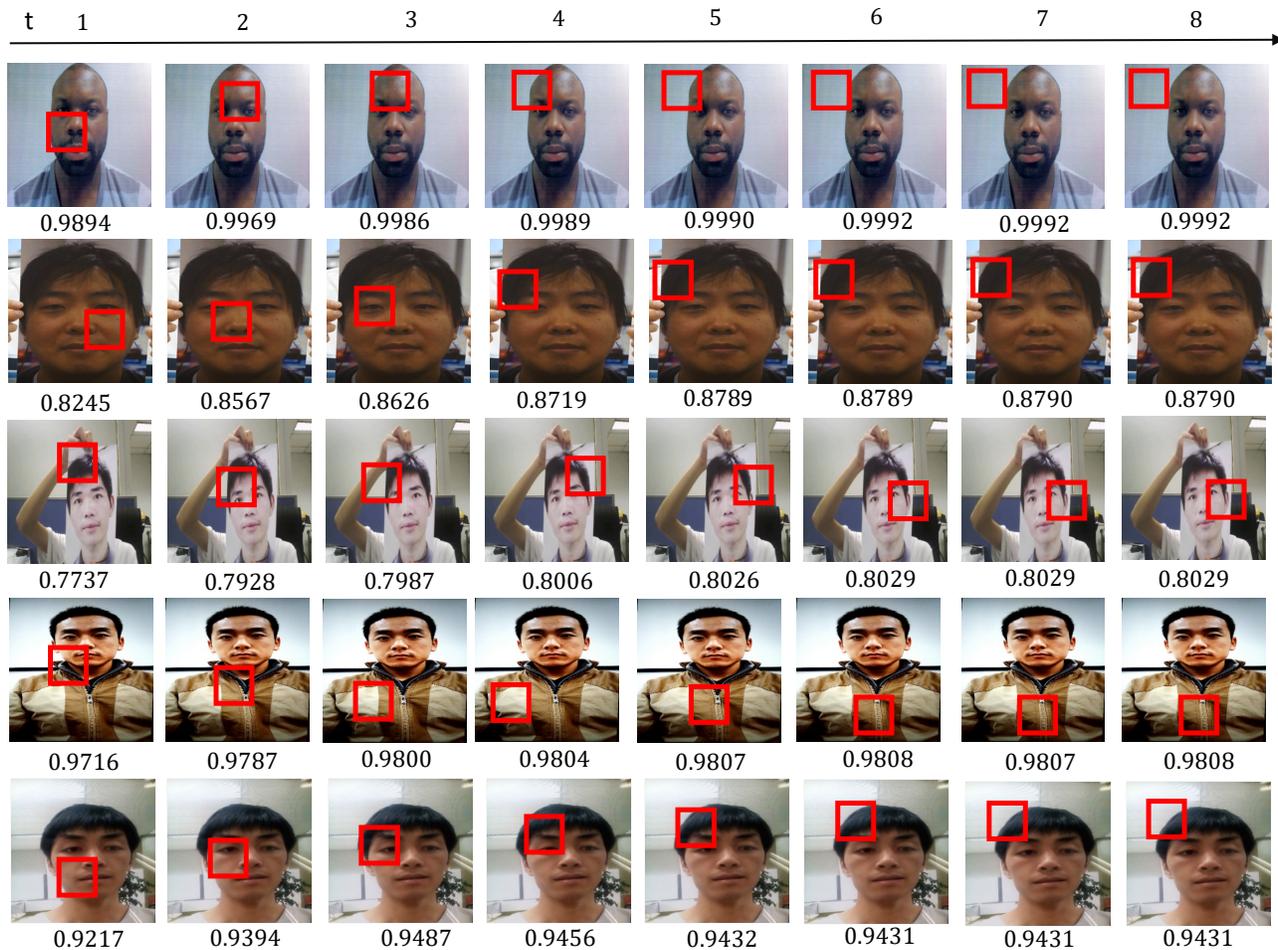


Fig. 8. Visualization of each step with the confidence score. The spoofing examples are correctly classified in the intra-database experiments. Each row represents an example, and each column shows the predicted location at a step. The index of an observation step is denoted by t . Under each image is the confidence score c_t . The first row is a replay attack example from the IDIAP REPLAY-ATTACK database. The second and third rows are printed paper attack example from the CASIA database. The fourth row is a replay video attack example from the OULU-NPU database. The fifth row is a replay video attack example from the ROSE-YOUTU database.



Fig. 9. Misclassification of spoofing examples. (a), (b), and (c) were misclassified as genuine faces in the intra-database experiments on the CAISA, ROSE-YOUTU, and OULU-NPU databases, respectively.

TABLE XVI
THE COMPUTATIONAL INFORMATION OF RESNET18 AND OUR METHOD.

Method	Params	MACs	FPS
ResNet18	11.18M	2.38G	150
Ours (T=4)	16.50M	2.42G	110
Ours (T=8)	16.50M	2.45G	70

V. CONCLUSION

We present a novel two-branch framework to explore spoofing clues for face anti-spoofing problem. The novelties of our work lie in two folds, 1) we propose to leverage CNN and RNN to extract both global and local information for the FAS task based on a single frame; 2) we propose a novel optimization strategy based on deep reinforcement learning, which is the first attempt in the FAS problem. We conduct extensive experiments on six different databases to evaluate our proposed framework. The experimental results on both intra- and cross-domain indicate that our proposed framework can generally achieve state-of-the-art performance compared with various state-of-the-art baselines, which demonstrate the effectiveness of our method.

REFERENCES

- [1] H. Li, W. Li, H. Cao, S. Wang, F. Huang, and A. C. Kot, “Unsupervised domain adaptation for face anti-spoofing,” *IEEE Transactions on Information Forensics and Security*, vol. 13, pp. 1794–1809, July 2018.

- [2] R. Nosaka, Y. Ohkawa, and K. Fukui, "Feature extraction based on co-occurrence of adjacent local binary patterns," in *Advances in Image and Video Technology* (Y.-S. Ho, ed.), (Berlin, Heidelberg), pp. 82–91, Springer Berlin Heidelberg, 2012.
- [3] I. Chingovska, A. Anjos, and S. Marcel, "On the Effectiveness of Local Binary Patterns in Face anti-Spoofing," in *Biometrics Special Interest Group*, pp. 1–7, 2012.
- [4] Z. Boulkenafet, J. Komulainen, and A. Hadid, "Face spoofing detection using colour texture analysis," *IEEE Transactions on Information Forensics and Security*, vol. 11, pp. 1818–1830, Aug 2016.
- [5] X. Tan, Y. Li, J. Liu, and L. Jiang, "Face liveness detection from a single image with sparse low rank bilinear discriminative model," in *European Conference on Computer Vision*, pp. 504–517, 2010.
- [6] D. Gragnaniello, G. Poggi, C. Sansone, and L. Verdoliva, "An investigation of local descriptors for biometric spoofing detection," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 4, pp. 849–863, 2015.
- [7] J. Yang, Z. Lei, and S. Z. Li, "Learn convolutional neural network for face anti-spoofing," *Computer Science*, vol. 9218, pp. 373–384, 2014.
- [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances in Neural Information Processing Systems*, pp. 1097–1105, 2012.
- [9] Y. Liu, A. Jourabloo, and X. Liu, "Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, (Salt Lake City, UT), pp. 389–398, 2018.
- [10] Z. Sun, L. Sun, and Q. Li, "Investigation in spatial-temporal domain for face spoof detection," pp. 1538–1542, 04 2018.
- [11] W. Wang, W. Luo, L. Bao, Y. Gao, D. Gong, S. Zheng, Z. Li, and A. Overwijk, "Face anti-spoofing: Model matters, so does data," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [12] K. Patel, H. Han, and A. K. Jain, "Secure face unlock: Spoof detection on smartphones," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 10, pp. 2268–2283, 2016.
- [13] J. Maatta, A. Hadid, and M. Pietikainen, "Face spoofing detection from single images using micro-texture analysis," in *2011 International Joint Conference on Biometrics (IJCB)*, pp. 1–7, IEEE, 2011.
- [14] J. Yang, Z. Lei, S. Liao, and S. Z. Li, "Face liveness detection with component dependent descriptor," in *2013 International Conference on Biometrics (ICB)*, pp. 1–6, IEEE, 2013.
- [15] Z. Boulkenafet, J. Komulainen, and A. Hadid, "Face Anti-Spoofing Using Speeded-Up Robust Features and Fisher Vector Encoding," *IEEE Signal Processing Letters*, vol. 24, no. 2, pp. 141–145, 2017.
- [16] J. Galbally and S. Marcel, "Face anti-spoofing based on general image quality assessment," in *2014 22nd International Conference on Pattern Recognition*, pp. 1173–1178, Aug 2014.
- [17] D. Wen, H. Han, and A. K. Jain, "Face spoof detection with image distortion analysis," *IEEE Transactions on Information Forensics and Security*, vol. 10, pp. 746–761, April 2015.
- [18] H. Li, S. Wang, and A. C. Kot, "Face spoofing detection with image quality regression," in *2016 Sixth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pp. 1–6, 2016.
- [19] A. Anjos and S. Marcel, "Counter-measures to photo attacks in face recognition: a public database and a baseline," in *2011 International Joint Conference on Biometrics (IJCB)*, pp. 1–7, IEEE, 2011.
- [20] J. Komulainen, A. Hadid, M. Pietikinen, A. Anjos, and S. Marcel, "Complementary countermeasures for detecting scenic face spoofing attacks," in *2013 International Conference on Biometrics (ICB)*, pp. 1–7, June 2013.
- [21] T. de Freitas Pereira, J. Komulainen, A. Anjos, J. M. De Martino, A. Hadid, M. Pietikäinen, and S. Marcel, "Face liveness detection using dynamic texture," *EURASIP Journal on Image and Video Processing*, vol. 2014, no. 1, p. 2, 2014.
- [22] L. Li, Z. Xia, A. Hadid, X. Jiang, H. Zhang, and X. Feng, "Replayed video attack detection based on motion blur analysis," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 9, pp. 2246–2261, 2019.
- [23] L. Feng, L.-M. Po, Y. Li, X. Xu, F. Yuan, T. C.-H. Cheung, and K.-W. Cheung, "Integration of image quality and motion cues for face anti-spoofing: A neural network approach," *Journal of Visual Communication and Image Representation*, vol. 38, pp. 451–460, 2016.
- [24] H. Chen, G. Hu, Z. Lei, Y. Chen, N. M. Robertson, and S. Z. Li, "Attention-based two-stream convolutional networks for face spoofing detection," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 578–593, 2020.
- [25] A. George, Z. Mostaani, D. Geissenbuhler, O. Nikisins, A. Anjos, and S. Marcel, "Biometric face presentation attack detection with multi-channel convolutional neural network," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 42–55, 2020.
- [26] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, pp. 1735–80, 12 1997.
- [27] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.
- [28] Z. Xu, S. Li, and W. Deng, "Learning temporal features using LSTM-CNN architecture for face anti-spoofing," in *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, pp. 141–145, 2016.
- [29] Y. Atoum, Y. Liu, A. Jourabloo, and X. Liu, "Face anti-spoofing using patch and depth-based CNNs," in *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pp. 319–328, Oct 2017.
- [30] A. Jourabloo, Y. Liu, and X. Liu, "Face de-spoofing: Anti-spoofing via noise modeling," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 290–306, September 2018.
- [31] H. Li, P. He, S. Wang, A. Rocha, X. Jiang, and A. C. Kot, "Learning generalized deep feature representation for face anti-spoofing," *IEEE Transactions on Information Forensics and Security*, vol. 13, pp. 2639–2652, Oct 2018.
- [32] R. Shao, X. Lan, J. Li, and P. C. Yuen, "Multi-adversarial discriminative deep domain generalization for face presentation attack detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 10023–10031, 2019.
- [33] Y. Liu, J. Stehouwer, A. Jourabloo, and X. Liu, "Deep tree learning for zero-shot face anti-spoofing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4680–4689, 2019.
- [34] H. Li, S. Wang, P. He, and A. Rocha, "Face anti-spoofing with deep neural network distillation," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 5, pp. 933–946, 2020.
- [35] Z. Li, H. Li, K. Lam, and A. C. Kot, "Unseen face presentation attack detection with hypersphere loss," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2852–2856, 2020.
- [36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [37] M. Lin, Q. Chen, and S. Yan, "Network in network," in *the Proceedings of International Conference on Learning Representations (ICLR)*, 2014.
- [38] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Advances in neural information processing systems*, pp. 1057–1063, 2000.
- [39] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3–4, pp. 229–256, 1992.
- [40] Z. Zhang, J. Yan, S. Liu, Z. Lei, D. Yi, and S. Z. Li, "A face anti-spoofing database with diverse attacks," in *IAPR International Conference on Biometrics*, pp. 26–31, 2012.
- [41] Z. Boulkenafet, J. Komulainen, L. Li, X. Feng, and A. Hadid, "OULU-NPU: A mobile face presentation attack database with real-world variations," in *IEEE International Conference on Automatic Face and Gesture Recognition*, May 2017.
- [42] I. Shah, *The Elephant in the Dark*. "Geneva University lectures, 1972/3," Octagon Press, 1974.
- [43] D. E. King, "Dlib-ml: A machine learning toolkit," *Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.
- [44] Z. Boulkenafet, J. Komulainen, Z. Akhtar, A. Benlamoudi, D. Samai, S. E. Bekhouche, A. Ouafi, F. Dornaika, A. Taleb-Ahmed, L. Qin, F. Peng, L. B. Zhang, M. Long, S. Bhilare, V. Kanhangad, A. Costa-Pazo, E. Vazquez-Fernandez, D. Perez-Cabo, J. J. Moreira-Perez, D. Gonzalez-Jimenez, A. Mohammadi, S. Bhattacharjee, S. Marcel, S. Volkova, Y. Tang, N. Abe, L. Li, X. Feng, Z. Xia, X. Jiang, S. Liu, R. Shao, P. C. Yuen, W. R. Almeida, F. Andalo, R. Padilha, G. Bertocco, W. Dias, J. Wainer, R. Torres, A. Rocha, M. A. Angeloni, G. Folego, A. Godoy, and A. Hadid, "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, Oct 2017.
- [45] T. de Freitas Pereira, A. Anjos, J. M. De Martino, and S. Marcel, "Can face anti-spoofing countermeasures work in a real world scenario?," in *2013 International Conference on Biometrics (ICB)*, pp. 1–8, IEEE, 2013.

- [46] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2921–2929, 2016.
- [47] IEEE Standards Committee, "IEEE Standard for Floating-Point Arithmetic," *IEEE Std 754-2008*, pp. 1–70, 2008.