# Learning One Class Representations for Face Presentation Attack Detection using Multi-channel Convolutional Neural Networks

Anjith George, *Member, IEEE* and Sébastien Marcel, *Senior Member, IEEE*

*Abstract*—Face recognition has evolved as a widely used biometric modality. However, its vulnerability against presentation attacks poses a significant security threat. Though presentation attack detection (PAD) methods try to address this issue, they often fail in generalizing to unseen attacks. In this work, we propose a new framework for PAD using a one-class classifier, where the representation used is learned with a Multi-Channel Convolutional Neural Network (*MCCNN*). A novel loss function is introduced, which forces the network to learn a compact embedding for *bonafide* class while being far from the representation of attacks. A one-class Gaussian Mixture Model is used on top of these embeddings for the PAD task. The proposed framework introduces a novel approach to learn a robust PAD system from *bonafide* and available (known) attack classes. This is particularly important as collecting *bonafide* data and simpler attacks are much easier than collecting a wide variety of expensive attacks. The proposed system is evaluated on the publicly available *WMCA* multi-channel face PAD database, which contains a wide variety of 2D and 3D attacks. Further, we have performed experiments with *MLFP* and *SiW-M* datasets using RGB channels only. Superior performance in unseen attack protocols shows the effectiveness of the proposed approach. Software, data, and protocols to reproduce the results are made available publicly.

*Index Terms*—Presentation Attack Detection, Convolutional Neural Network, Face Recognition, Anti-spoofing, Reproducible Research, Unseen Attack Detection.

## I. INTRODUCTION

**F**ACE recognition has proved to be a beneficial modality for biometric authentication. One of the main reasons for the widespread use of face recognition systems is its non-intrusive nature of acquisition and ease of use [1]. Face recognition systems have matured a lot in recent years, and several approaches have reported human parity in the identification rate in 'in the wild' conditions [2]. However, a critical security issue undermining the widespread use of face recognition technology is its vulnerability to presentation attacks (a.k.a spoofing attacks) [3], [4].

Presentation attack refers to an attack using an instrument with the intention to affect the normal operation of the biometric system. Often, features such as color, texture [5], [6], motion [7], and physiological cues [8], [9] and CNN based methods [10] are used for detection of attacks like 2D prints and replays. However, detection of sophisticated attacks like 3D masks and partial attacks are challenging and poses a serious threat to the reliability of face recognition systems.

A. George and S. Marcel are in Idiap Research Institute, Centre du Parc, Rue Marconi 19, CH - 1920, Martigny, Switzerland. (e-mail: {anjith.george,sebastien.marcel}@idiap.ch)
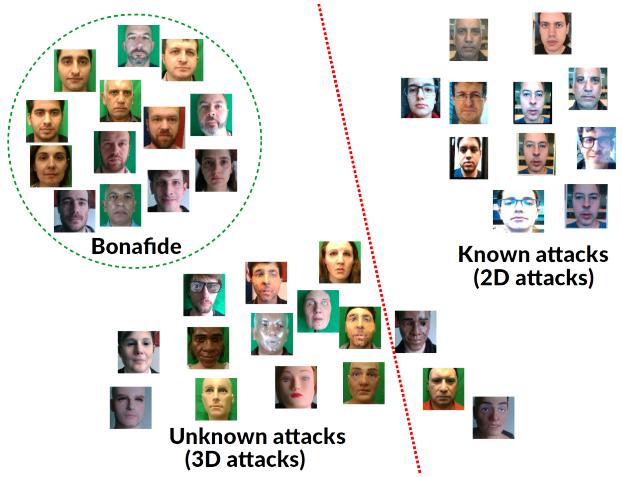


Fig. 1. Illustration of the embedding space with known and unknown attack classes. The red dotted line shows the decision boundary learned when only *bonafide* and known attacks are present in the training set, this results in misclassification of unknown attacks as *bonafide*. If a one class decision boundary (green-dotted lines) is learned, then both known and unknown attacks can be classified correctly.

Most of the presentation attack detection (PAD) methods available in prevailing literature try to solve the problem for a limited number of presentation attack instruments and on visible spectrum images [3]. Though some success has been achieved in addressing 2D presentation attacks, performance of the algorithms in realistic 3D masks and other kinds of attacks is poor. With the increase in quality of attack instruments, it becomes harder to discriminate between *bonafide* and PAs in the visible spectrum alone. Moreover, considering a real-world situation with a wide variety of 2D, 3D, and partial attacks, PAD in visual spectra alone is challenging and inadequate for security-critical applications. Partial attacks refer to attacks where the attack instrument covers only a part of the face. These attacks are much harder to detect as they appear similar to *bonafide* in most of the face regions, and they can fool holistic liveliness detection systems easily. Multi-channel methods have been proposed as an alternative [11], [12], [13], [14], [15], [16], [15], [17], since they use complementary information from different channels to improve the discrimination between *bonafide* and attacks. In the multi-channel scenario, the additional channels used can be any modality which can provide complementary representation such as depth, infrared, and thermal channels. Multi-channel PAD approaches are more promising in the context of a wide variety of attacks since they make PAD systems harder to fool.

Even with the use of multiple channels, one of the main issues with PAD is its poor generalization to unseen attacks [14]. This is particularly important, since at the time of developing a PAD system, anticipating all possible attacks is impossible. Malicious attackers can always come up with new attacks to fool the PAD systems. In such situations, PAD systems which are robust against unseen attacks are of paramount importance. Moreover, while it is comparatively easy to collect data for attacks like 2D prints and replays, making replicas of challenging presentation attack instruments (PAI) like silicone mask are often very costly [18] and resource-intensive. In this context, it will be ideal to have a framework which can be trained with *bonafide* alone, or with a combination of *bonafide* and easy to manufacture PAIs.

In real-world scenarios, it can be assumed that all presentation attacks are unseen, as it is not possible to foretell all the variations a PAD system could encounter a priori. A toy example of the decision boundary in an unseen attack scenario is illustrated in Fig. 1. Performances in typical PAD databases may not be representative of the performance of a PAD system in real-world conditions. This necessitates the PAD algorithms to be robust against unseen attacks. Since it is easy (in effort and cost) to collect data from more straightforward attacks compared to complex PAIs, we try to learn the representation leveraging the information from PA classes which are available at the training stage (while not over-fitting on the available attacks). To achieves this, we propose a one-class classifier based framework, where the feature representation is learned with a CNN to have discriminative properties. The core of the framework is a multi-channel CNN trained to learn the embedding using a specific loss function. The proposed approach aims at learning a compact representation for the *bonafide* class while leveraging the discriminative information for PAD task.

The main contributions of the paper are listed below.

- A novel multi-channel one-class classifier-based approach is proposed for unseen attack detection.
- A novel loss function is proposed which learns a compact and discriminative representation of the face for PAD task, leveraging the information provided from known attacks.

The features used in the one class classifier are learned with a multi-channel CNN framework. The proposed approach was evaluated in *known* and *unseen* attack protocols in *WMCA* database containing a wide variety of 2D and 3D attacks, and performed significantly better than baselines in *unseen* protocols. We have also performed experiments using RGB channel in *MLFP* and *SiW-M* datasets.

Additionally, the source code and protocols to reproduce the results are made available publicly and are accessible at the following link [1].

The rest of the paper is organized as follows. Section 2 describes the related work with a particular focus on unseen attack detection. Section 3 outlines the proposed framework. Extensive evaluations, comparison with baseline methods, and ablation studies are shown in section 4. Section 5 discusses

---

[1] Source code: https://gitlab.idiap.ch/bob/bob.paper.oneclass_mccnn_2019

the importance of the results, and Section 6 presents the conclusions.

## II. RELATED WORK

Majority of the literature in face PAD is mainly focused on 2D attacks and uses feature-based methods [5], [6], [7],[8], [9] or CNN based methods. Recently, CNN based methods have been more successful as compared to feature-based methods [19], [10], [20], [21]. These methods usually leverage the quality degradation during 'recapture' and are often useful only for the detection of attacks like 2D prints and replays. Sophisticated attacks like 3D masks are more challenging and pose serious threat to the reliability of face recognition systems.

Most of these methods handle the PAD problem as binary classification, which results in classifiers over-fitting to the known attacks resulting in poor generalization to unseen attacks. We focus the further discussion on the detection of unseen attacks. However, it is imperative that methods working for unseen attacks must perform accurately for known attacks as well. One naive solution for such a task is one-class classifiers (OCC). OCC provides a straightforward way of handling the unseen attack scenario by modeling the distribution of the *bonafide* class alone.

Arashloo *et al*.[22] and Nikisins *et al*. [23] have shown the effectiveness of one class methods against unseen attacks. Even though these methods performed better than binary classifiers in an unseen attack scenario, the performance in known attack protocols was inferior to that of binary classifiers. Xiong *et al*. [24] proposed unseen PAD methods using auto-encoders and one class classifiers with texture features extracted from images. However, the performance of the methods compared to recent CNN based methods is very poor. CNN based methods outperform most of the feature-based baselines for PAD task. Hence there is a clear need of one class classifiers or anomaly detectors in the CNN framework. One of the drawbacks of one class model is that they do not use the information provided by the known attacks. An anomaly detector framework which utilizes the information from the known attacks could be more efficient.

Perera and Patel [25] presented an approach for one-class transfer learning in which labelled data from an unrelated task is used for feature learning. They used two loss functions, namely descriptive loss, and compactness loss to learn the representations. The data from the class of interest is used to calculate the compactness loss whereas an external multi-class dataset is used to compute the descriptive loss. Accuracy of the learned model in classification using another database is used as the descriptive loss. However, in the face PAD problem, this approach would be challenging since the *bonafide* and attack classes appear very similar.

Fatemifar *et al*. [26] proposed an approach to ensemble multiple one-class classifiers for improving the generalization of PAD. They introduced a class-specific normalization scheme for the one class scores before fusion. Seven regions, three one class classifiers and representations from three CNNs were used in the pool of classifiers. Though their method

achieved better performance as compared to client independent thresholds, the performance is inferior to CNN based state of the art methods. Specifically, many CNN based approaches have achieved 0% HTER in Replay-Attack and Replay-Mobile datasets. Moreover, the challenging unseen attack scenario is not evaluated in this work.

Pérez-Cabo *et al*. [27] proposed a PAD formulation from an anomaly detection perspective. A deep metric learning model is proposed, where a triplet focal loss is used as a regularization for 'metric-softmax', which forces the network to learn discriminative features. The features learned in such a way is used together with an SVM with RBF kernel for classification. They have performed several experiments on an aggregated RGB only datasets showing the improvement made by their proposed approach. However, the analysis is mostly limited to RGB only models and 2D attacks. Challenging 3D and partial attacks are not considered in this work. Specifically, the effectiveness in challenging unknown attacks (2D vs 3D) is not evaluated.

Recently, Liu *et al*. [28] proposed an approach for the detection of unknown spoof attacks as Zero-Shot Face Anti-spoofing (ZSFA). They proposed a Deep Tree Network (DTN) which partitions the attack samples into semantic sub-groups in an unsupervised manner. Each tree node in their network consists of a Convolutional Residual Unit (CRU) and a Tree Routing Unit (TRU). The objective is to route the unknown attacks to the most proper leaf node for correctly classifying it. They have considered a wide variety of attacks in their approach and their approach achieved superior performance compared to the considered baselines.

Jaiswal *et al*. [29] proposed an end to end deep learning model for PAD which used unsupervised adversarial invariance. In their method, the discriminative information and nuisance factors are disentangled in an adversarial setting. They showed that by retaining only discriminative information, the PAD performance improved for the same base architecture. Mehta *et al*. [30] trained an Alexnet model with a combination of cross-entropy and focal losses. They extracted the features from Alexnet and trained a two-class SVM for PAD task. However, results in challenging datasets such as OULU and SiW were not reported.

Recently Joshua and Jain [31] utilized multiple GANs for spoof detection in fingerprints. Their method essentially consisted of training a DCGAN [32] using only the *bonafide* samples. At the end of the training, the generator is discarded, and the discriminator is used as the PAD classifier. They combined the results from different GANs operating on different features. However, this approach may not work well for face images as the recaptured images look very similar to the *bonafide* samples.

In safety critical applications, extended range methods have been proposed over the years [11], [33], [12], [34], [35], [18], [13], [14] to achieve reliable PAD performance. Even these methods fail in generalizing to unseen attacks.

Wang *et al*. [36] proposed multimodal face presentation attack detection with a ResNet based network using both spatial and channel attentions. Specifically, the approach was tailored for the *CASIA-SURF* [37] database which contained RGB, near-infrared and depth channels. The proposed model is a multi-branch model where the individual channels and fused data are used as inputs. Each input channel has its own feature extraction module and the features extracted are concatenated in a late fusion strategy. Followed by more layers to learn a discriminative representation for PAD. The network training is supervised by both center loss and softmax loss. One key point is the use of spatial and channel attention to fully utilize complementary information from different channels. Though the proposed approach achieved good results in the *CASIA-SURF* database, the challenging problem of unseen attack detection is not addressed.

Parkin *et al*. [38] proposed a multi-channel face PAD network based on ResNet. Essentially, their method consists of different ResNet blocks for each channel followed by fusion. Squeeze and excitation modules (SE) are used before fusing the channels, followed by remaining residual blocks. Further, they add aggregation blocks at multiple levels to leverage inter-channel correlations. Their approach achieved state of the art results in *CASIA-SURF* [37] database. However, the final model presented in is a combination of 24 neural networks trained with different attack specific folds, pre-trained models and random seeds, which would increase the computation greatly.

From the discussions above, it can be seen that one class classifiers could be a good alternative for binary classification in PAD task. However, the features used for one class classifiers should be discriminative and compact to outperform binary classification.

## III. PROPOSED METHOD

From a practical viewpoint, it is not possible to anticipate all the possible types of attacks and to have them in the training set. This, in turn, make the PAD task an unseen classification problem in a broad sense. In general, we can even consider attacks coming from different replay devices as unseen attacks. Typically, one class classifiers are well suited for such outlier detection tasks. However, in practice, the performance of one class classifiers are inferior compared to binary classifiers for known attacks, since they do not leverage useful information from the known attacks. Ideally, the PAD system should perform well in both known and unseen attack scenarios.

Clearly, there is a necessity of a method which can learn a compact one class representation while utilizing the discriminative information from known attacks. While the collection of attacks could be difficult and costly, collecting *bonafide* samples are rather easy. A new classification strategy is required to handle the realistic scenario where a limited variety of attack classes are available.

Though one class classifiers (*OCC*) offers a way to model the *bonafide* class, the efficient use of *OCC* requires the feature representation to be compact while containing discriminative information for PAD task. In the proposed framework, we use a CNN based approach to learn the feature representation. A novel loss function is proposed to learn a representation of *bonafide* samples leveraging the known attack classes.

## A. Formulation of One Class Contrastive Loss (OCCL)

Consider a typical CNN architecture for PAD, where the output layer contains one node and the loss function used is Binary Cross Entropy (*BCE*), which is defined as:

$$\mathcal{L}_{BCE} = -(y \log(p) + (1-y) \log(1-p)) \qquad (1)$$

where $y$ is the ground truth, ($y = 0$ for attack and $y = 1$ for *bonafide*) and $p$ is the probability.

When trained only with *BCE* loss, the network learns a decision boundary based on the *bonafide* and attacks present in the training set. However, it may not generalize when encountered with an unseen attack in the test time as it could be over-fitted to attacks which are 'known' from the training set.

To overcome this issue, we propose the 'One-Class Contrastive Loss' (*OCCL*) function which operates on the embedding layer. Proposed One-Class Contrastive Loss (*OCCL*) function is used as an auxiliary loss function in conjunction with binary cross-entropy loss. The feature map obtained from the penultimate layer of the CNN is used as the embedding. The loss function is inspired from center-loss [39] and contrastive loss [40], which are usually used in the face recognition applications.
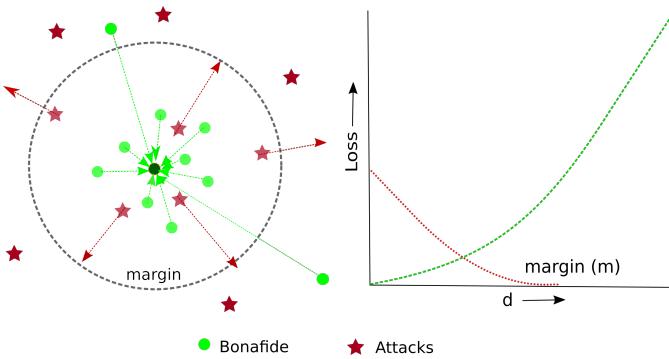


Fig. 2. Loss functions acting on the embedding space, left) *bonafide* representations are pulled closer to the center of *bonafide* class (green), while the attack embeddings(red) are forced to be beyond the margin. The attack samples outside the margin does not contribute to the loss, right) The loss as a function of distance from the *bonafide* center.

In face recognition applications, center loss is used as an additional auxiliary loss function, the task of the center loss is to minimize the distance of the embeddings from their corresponding class centers. The center loss is defined as:

$$\mathcal{L}_{center} = \frac{1}{2} \sum_{i=1}^{m} \|x_i - c_{y_i}\|_2^2 \qquad (2)$$

Where $L_{center}$ denotes the center loss, $m$ the number of training samples in a mini-batch, $x_i \in R_d$ denotes the $i^{th}$ training sample, $y_i$ denotes the label, and $c_{y_i}$ denotes the $y_i^{th}$ class center in the embedding space.

The main issue with center loss in the PAD application is that the loss function penalizes for large intra-class distances and does not care about the inter-class distances. Contrastive center loss [41] tries to solve this issue by adding the distance between classes (inter-class) in the formulation. However, for

the PAD problem, modeling the attack class as a cluster and finding a center for the attack class is not trivial. The attacks could be of different categories: 2D, 3D, and partial attacks, and it is not ideal forcing them to cluster together in the embedding space. It is only necessary to have the embeddings of attacks far from *bonafide* cluster in the embedding space. Hence, we put the compactness constraint only on the *bonafide* class, while forcing the embeddings of PAs to be far from that of *bonafide*.

To formulate the loss function, we start with the equation for contrastive loss function proposed by Lecun *et al.* [40].

$$\mathcal{L}_{Contrastive}(W, Y, X^1, X^2) = (1-Y)\frac{1}{2}D_W^2$$
$$+ Y\frac{1}{2}max(0, m - D_W)^2 \qquad (3)$$

Where $W$ is the network weights, $X^1, X^2$ are the pairs and $Y$ the label of the pair, i.e., whether they belong to the same class or not. $m$ is the margin, and $D_W$ is the distance function between two samples. The data is provided as pairs $(X^1, X^2)$ and the distance function $D_W$ can be computed as the Euclidean distance.

$$D_W = \sqrt{\|X^1 - X^2\|_2^2} \qquad (4)$$

Now, in our loss formulation, the critical difference is how we define $D_W$. In the original contrastive loss, $D_W$ is the distance between samples. In our case, we need the representation of *bonafide* samples to be compact in an embedding space. At the same time, we want to maximize the distance between *bonafide* cluster and attack samples in the embedding space. This can be achieved by defining $DC_W$ to be the distance from the center of *bonafide* class as follows.

$$DC_W = \sqrt{\|X^i - c_{BF}\|_2^2} \qquad (5)$$

Where $X^i$ is the embedding for $i^{th}$ sample, and $c_{BF}$ is the center of *bonafide* class in the embedding space.

The center of the *bonafide* class is updated in every mini-batch during training as follows.

$$c_{BF} = \hat{c}_{BF}(1-\alpha) + \alpha\frac{1}{N}\sum_{i=1}^{N} e_i \qquad (6)$$

Where $c_{BF}$ and $\hat{c}_{BF}$ denotes the new and old *bonafide*-centers. $\alpha$ is a scalar which prevents sudden changes in the class centers in mini-batch. $e_i$ denotes the difference between embeddings for the *bonafide* samples in the current mini-batch compared to the previous center, and $N$ denotes the number of *bonafide* samples in the mini-batch.

Combining the equations, our auxiliary loss function becomes:

$$\mathcal{L}_{OCCL}(W, Y, X) = Y\frac{1}{2}DC_W^2$$
$$+ (1-Y)\frac{1}{2}max(0, m - DC_W)^2 \qquad (7)$$

Where $DC_W$ denotes the Euclidean distance between the samples and the *bonafide* class center, $Y$ denotes the ground
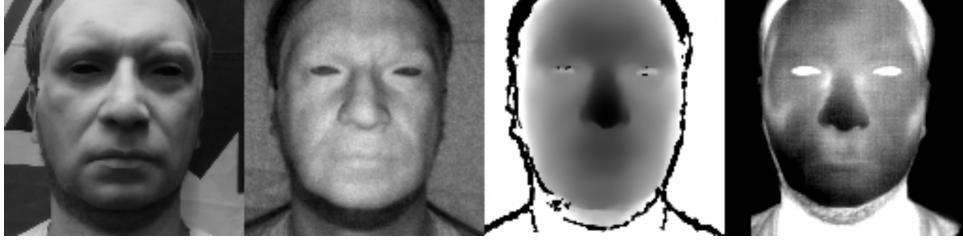
Fig. 3. Preprocessed images from a rigid mask attack; channels showed are gray-scale, infrared, depth, and thermal, respectively. Channels were preprocessed with face detection, alignment and normalization.

truth, i.e., $Y = 0$ for attacks and $Y = 1$ for *bonafide* (note the change in labels from the standard notation due to the ground truth convention). It is to be noted that, the proposed loss function ***does not*** require pairs of samples, which is a requirement in usage of contrastive loss. This makes it easier to train the model without requiring an explicit selection of pairs during training.

This auxiliary loss makes the representation of *bonafide* compact pushing it closer to the center of *bonafide* class and penalizes attack samples which are closer than the margin $m$. Attack samples which are farther than the margin $m$ are not penalized. An illustration of the loss functions acting on the embeddings of *bonafide* and attack samples are shown in Fig. 2.

We combine the proposed loss function with standard binary cross entropy for training. The combined loss function to minimize is given as:

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_{BCE} + \lambda\mathcal{L}_{OCCL} \qquad (8)$$

Where $\mathcal{L}$ denotes the total loss for the CNN. $\mathcal{L}_{BCE}$ and $\mathcal{L}_{OCCL}$ denotes the binary cross entropy, and one-class contrastive loss respectively. $\lambda$ denotes a scalar value to set the weight for each loss functions. In our experiments we set the value of $\lambda$ as 0.5.

The combined loss function $\mathcal{L}$ tries to learn a decision boundary between the available attacks and *bonafide* while the auxiliary loss tries to make the feature representation of the *bonafide* compact in the embedding space. We expect the decision boundary learned in this fashion to be more robust in unseen attacks compared to the network learned only with *BCE*. The embedding obtained in this manner is used with a one-class classifier for the PAD task.

### B. Components of the proposed framework

Different stages of the proposed framework are described below.

*1) Preprocessing:* Before using the data from the sensors, a preprocessing stage consisting of face detection, alignment, and normalization is performed. MTCNN algorithm [42] was used for face detection in the color channel followed by face landmark detection using Supervised Descent Method (SDM) [43]. After these stages, the face image is aligned and converted to gray-scale with a resolution of $128 \times 128$ pixels. Since all the channels are aligned, these face locations are utilized for the alignment of other non-RGB channels as well.

Also, normalization using Mean Absolute Deviation (MAD) [44] is performed to convert the raw 16-bit values to the 8-bit range. An example image after preprocessing stage is shown in Fig. 3.

*2) Network architecture and training:* Since the data used is multi-channel, we use a multi-channel PAD framework called 'Multi-Channel Convolutional Neural Network'(*MCCNN*) proposed in [14] as our base network. The main idea in *MCCNN* was to use the joint representation from multiple channels for PAD task, leveraging a pretrained face recognition network. The *MCCNN* architecture constituted of an extended version of LightCNN model [45] adapted specifically for multi-channel PAD task. A pretrained LightCNN face recognition model was extended to accept multiple channels, and the embeddings from all channels were concatenated, and two fully connected layers were added on top of this joint representation layer for PAD task. The advantage in this architecture is that only lower layer features (which are known as Domain Specific Units (DSU) [46] ) and higher-level fully connected layers are adapted in the training phase. The first fully connected layer contains ten nodes, and the second layer contains only one output node. The higher-level features in the LightCNN part are shared among all the modalities. This approach has two main advantages; first, there is a smaller number of parameters since the high-level features are shared across modalities, second, adapting only DSUs and final fully connected layers reduce possible over-fitting since PAD databases are typically small in size. An optimal set of layers to be adapted was obtained empirically and was used in the baseline *MCCNN* and the proposed approach.

In our proposed approach, we use the same *MCCNN* architecture, and the output from the penultimate fully connected layer was used as the embeddings. To quantify the effectiveness of our approach, we perform experiments on the *MCCNN* architecture, while using both embeddings and the final output for the loss computation. An illustration of the proposed framework is shown in Fig. 4. At the time of training, both losses are used, and the model corresponding to the lowest validation score is selected. It is to be noted that, at the time of CNN training, both *bonafide* and (known) attack samples are used. After the CNN training, the network weights are frozen, and the *bonafide* samples are feed-forwarded to obtain the embeddings.

*3) One-Class Gaussian Mixture Model:* After the training of *MCCNN* with *BCE* and *OCCL*, the trained weights of the network are frozen, and it is used as a fixed feature
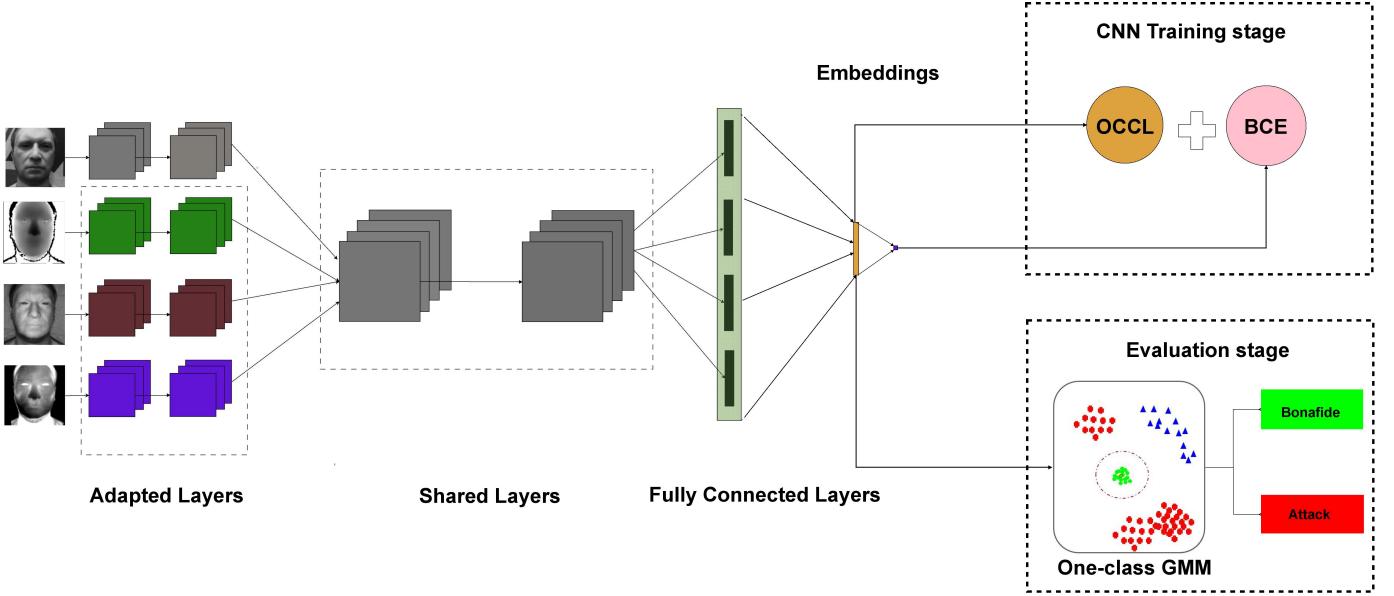
Fig. 4. Schematic diagram of the proposed framework. The CNN architecture is trained with two losses and then used as a fixed feature extractor with frozen weights. The one-class GMM is trained using the embeddings obtained from *bonafide* class alone.

extractor for the PAD task. Now that a compact representation is available, the objective is to learn a one-class classifier using the features obtained. We use One-Class Gaussian Mixture Model for this task. The one class GMM is a generative approach which is used for modeling the distribution of the *bonafide* class in the proposed framework.

A Gaussian Mixture Model is defined as the weighted sum of $K$ multivariate Gaussian distributions as:

$$p(x|\Theta) = \sum_{k=1}^{K} w_k \mathcal{N}(x; \mu_k, \Sigma_k), \qquad (9)$$

where $\Theta = \{w_k, \mu_k, \sigma_k\}_{\{k=1,...,K\}}$ are the weights, means and the covariance matrix of the GMM.

Expectation-Maximization (EM) [47] was used to compute the parameters of the GMM. A full covariance matrix is computed for each component, and the number of components to use was empirically selected as five ($K = 5$).

During the training phase, embeddings obtained from *bonafide* class only are used to train the One-Class GMM.

In test time, a sample is first forwarded though the network to obtain the embedding $x$, and then fed to the One-Class GMM to obtain the log-likelihood score as follows:

$$score = log(p(x|\Theta)) \qquad (10)$$

In summary, the proposed framework can be considered as a one-class classifier based framework for PAD. The crucial distinction is that, the features used are ***learned***. The loss function proposed forces the CNN to learn a compact representation for the *bonafide* class leveraging the information from known attack classes. The algorithm for training the framework is shown in Algorithm 1.

---

**Algorithm 1:** Algorithm for training the proposed framework

**Data:** $(x_i, y_i)$, where $x_i$ is multi-channel input and $y_i \in 0, 1$; 0 – for attack and 1– for *bonafide*

**Result:** $W_C$ – CNN weights, $\Theta_{GMM}$ – Parameters of GMM

1 **Constants** : $\lambda$ – weighting factor, $\mu$ – learning rate
2 **Initialize** : $C_{BF}$ – center of *bonafide* class, $W_C$ – initial weights of CNN from pretrained model
3 **for** *mini-batch* ← *1 to P* **do**
4      Forward $x_i$ through the CNN
5      Compute the combined loss:
     $\mathcal{L} = (1 - \lambda)\mathcal{L}_{BCE} + \lambda\mathcal{L}_{OCCL}$
6      Back-propagate the loss and update the weights of DSUs and FC layers
7      Update the *bonafide* center:
8      $c_{BF} = \hat{c}_{BF}(1 - \alpha) + \alpha\frac{1}{N}\sum_{i=1}^{N} e_i$
9 **end**
10 Forward $x_j$ (*bonafide*, where $y_j = 1$) through the CNN to obtain Embeddings $E_j$
11 Estimate parameters of GMM from $E_j$:
12 $\Theta_{GMM}= (w_k, \mu_k, \Sigma_k)$
13 **Parameters**← $(W_C, \Theta_{GMM})$

---

### C. Implementation details

To increase the number of samples, data augmentation using random horizontal flips with a probability of 0.5 was used in training. Adam Optimizer [48] was used to minimize the combined loss function. Learning rate of $1\times10^{-4}$ and a weight decay parameter of $1 \times 10^{-5}$ was used. The network was trained for 50 epochs on GPU grid with a batch size of 32. The model corresponding to minimum validation loss in the *dev* set is selected as the best model. For the four-channel models, the MCCNN architecture has about 13.1M parameters and about

14.5 GFLOPS. The implementation was done using PyTorch [49] library.

## IV. EXPERIMENTS

In order to evaluate the effectiveness of the proposed approach, we have performed experiments in three publicly available databases, namely *WMCA* [14], *MLFP* [35], and *SiW-M* [28] datasets. Recently published *CASIA-SURF* [37] database also consists of multi-channel data, namely color, depth, and infrared channels with a limited set of attack instruments. However, the raw data from the sensors were not publicly available; in the publicly available version of the database, images were masked and scaled with custom preprocessing reducing the dynamic range of depth and infrared channels severely. Moreover, there was no guaranteed alignment between the channels. Therefore we can't use our framework with *CASIA-SURF* database due to the mentioned limitations.

### A. WMCA dataset

We have conducted an extensive set of experiments on *Wide Multi-Channel presentation Attack* (*WMCA*) [2] database, which contains a total of *1679* video samples of *bonafide* and attack attempts from *72* identities. The database contains information from four different channels collected simultaneously, namely, color, depth, infrared, and thermal channels. The data was collected using two consumer devices, Intel® RealSense™ SR300 capturing RGB-NIR-Depth streams, and Seek Thermal CompactPRO for the thermal channel. The database contained around eighty different PAIs constituting seven different categories of attacks: print, replay, funny eyeglasses, fake head, rigid mask, flexible silicone mask, and paper masks. The RGB visualization of the attack categories is shown in Fig. 5 and the different sessions in Fig. 6. Detailed information about the *WMCA* database can be found in the publication [14]. The statistics of the number of samples in each category and their types are shown in Table I. We have made challenging protocols in the *WMCA* dataset to perform an extensive set of evaluations emulating real-world unseen attack scenarios.

*1) Protocols in WMCA:* To test the performance of the algorithm in known and unseen attack scenarios, we created three protocols in the *WMCA* dataset. The protocols are described below.

- **grandtest** : This is the exact same *grandtest* protocol available with *WMCA* database, here all the attack types are present in almost equal proportions in the *train*, *development* and *evaluation* sets. The attack types and *bonafide* samples are divided into three folds, and the client ids are disjoint across the three sets. Each presentation attack instrument had a separate client id. The train, dev, eval splits were made in such a way that a specific PA instrument will appear in only one fold.
- **unseen-2D** : In this protocol, we use same splits as *grandtest* and removed all 2D attacks from *train* and *development* groups. *Evaluation* set contains only *bonafide*

[2]Database available at : https://www.idiap.ch/dataset/wmca



Fig. 5. Attack categories in *WMCA* dataset, only RGB images are shown. Print and Replay constitutes the 2D attacks and all others are 3D attacks (Image taken from [14]).



Fig. 6. Different sessions in *WMCA* dataset, only RGB images are shown. A total of six sessions was used the *WMCA* (Image taken from [14])

and 2D attacks. This emulates the performance of a system when encountered with 2D attacks which was not seen in training.

- **unseen-3D** : In this protocol, we use same splits as *grandtest* and removed all 3D attacks from *train* and *development* groups. *Evaluation* set contains only *bonafide* and 3D attacks. This emulates the performance of a system when encountered with 3D attacks which were not seen in training. This is the most challenging protocol as the model sees only the simpler 2D attacks in training and encounter challenging 3D attacks in testing.

While the *grandtest* protocol emulates the known attack scenario, other protocols emulate the unseen attack scenario. All protocols are made available publicly.

TABLE I
STATISTICS OF ATTACKS IN *WMCA* DATABASE

| PA Category | Type | #Presentations |
|---|---|---|
| *bonafide* | - | 347 |
| glasses | Partial | 75 |
| print | 2D | 200 |
| replay | 2D | 348 |
| fake head | 3D | 122 |
| rigid mask | 3D | 137 |
| flexible mask | 3D | 379 |
| paper mask | 3D | 71 |
| **TOTAL** | | **1679** |

### B. MLFP dataset

MLFP dataset [35] consists of attacks captured with seven 3D latex masks and three 2D print attacks. The dataset contains videos captured from color, thermal and infrared channels. Since channels were captured individually in different recording sessions, multi-channel approaches are not trivial. Also, the alignment of channels is not possible since they are not collected simultaneously. Hence, we only use the RGB videos from the MLFP dataset for our experiments. The database contains videos of 10 subjects wearing both print and latex masks. There are 440 videos are consisting of both attacks and *bonafide* for the RGB channel.

*1) Protocols in MLFP:* To emulate known and unseen attack scenarios, we created three new protocols in the *MLFP* dataset. There are two types of attacks, namely print and mask. Only two sets, i.e., *train* and *evaluation* are created due to the small size of the dataset. We used a subset of the train set (10%) for model selection. The protocols are described below.

- **grandtest** : This protocol emulates the known attack scenario. Both the attacks are present in both *train* and *evaluation* set. However, the subjects and the PAs are disjoint across the two sets.
- **unseen-print** : In this protocol, only *bonafide* and mask attacks are present in *train* set; the *evaluation* set contains only *bonafide* and print attacks. This emulates unseen attack scenario.
- **unseen-mask** : In this protocol, only *bonafide* and print attacks are present in *train* set; the *evaluation* set contains only *bonafide* and mask attacks. This protocol also emulates unseen attack scenario.

### C. SiW-M dataset

The Spoof in the Wild database with Multiple Attack Types (*SiW-M*) [28] consists of a wide variety of attacks captured only in RGB spectra. The database consists of images from 493 subjects, and a total of 660 *bonafide* and 968 attack samples. A total of 1628 files, consisting of 13 different attack types, collected in different sessions, pose, lighting, and expression (PIE) variations. The attacks consist of various types of masks, makeups, partial attacks, and 2D attacks. The videos are available in 1080P resolution.

*1) Protocols in SiW-M:* To emulate unseen attack scenarios, we use the leave-one-out (LOO) testing protocols available

with the *SiW-M* [28] dataset. The protocols consists of only *train* and *eval* sets. In each LOO protocol, the training set consists of 80% percentage of the live data and 12 types of spoof attacks. The evaluation set consists of 20% of *bonafide* data and the attack which was left out in the training phase. The subjects in *bonafide* sets are disjoint in *train* and *evaluation* sets. A subset of the train set (5%) was used for model selection. Additionally, we have created a *grandtest* protocol, specifically for cross-database testing which contains all the attack types in all the folds.

### D. Evaluation metrics

We report the standardized ISO/IEC 30107-3 metrics [4], Attack Presentation Classification Error Rate (APCER), and Bonafide Presentation Classification Error Rate (BPCER), and Average Classification Error Rate (ACER) in the $test$ set. A BPCER threshold of 1% is used for computing the threshold in $dev$ set. The APCER and BPCER in both *dev* and *eval* sets are also reported. Additionally, the ROC curves for experiments are also shown in all the protocols. For the *MLFP* dataset, we report only EER in the *evaluation* set since only two sets are available. For SiW-M database, we apply a threshold selected a-priori in all protocols, for computing the metrics, to be comparable with the results in [28].

### E. Baselines

We have implemented three feature-based baselines and two CNN based baselines. For a fair comparison, all the benchmarks are multi-channel methods and use the same four channels. Besides, an RGB only CNN model is also added for comparison. A short description of the baselines along with the acronyms used are shown below:

- *MC-RDWT-Haralick-SVM*: This baseline is the multi-channel extension of the RDWT-Haralick-SVM approach proposed in [35]; the images from all channels are stacked together after preprocessing. For each channel, the image is divided into a $4 \times 4$ grid, and Haralick [50] features obtained from the RDWT decompositions are concatenated from all the grids in all channels to get the joint feature vector. The joint feature is used with a linear SVM for PAD.
- *MC-RDWT-Haralick-GMM*: Here, the feature extraction stage is same as *MC-RDWT-Haralick-SVM*; however, the classifier used is one class GMM. Only *bonafide* samples are used in training this model. This model is added to show the performance of one class models in unseen attack scenarios.
- *MC-LBP-SVM*: Here, again, the same preprocessing is performed on all the channels first. After this, Spatially enhanced histograms of LBP representation from all the component channels are computed and concatenated to a feature vector. The features extracted are fed to an SVM for PAD task.
- *DeepPixBiS* : This is a CNN based system [10] trained using both binary and pixel-wise binary loss function. This model only uses RGB information for PAD.

- *MC-ResNetPAD*: We reimplemented the architecture from [38] extending it to four channels, based on their open-source implementation [3]. This approach obtained the first place solution in the 'CASIA-SURF' challenge. For a fair comparison, instead of using an ensemble we used the best pretrained model as suggested in [38].
- *MCCNN(BCE)* : This is the multi-channel CNN system described in [14], which achieved state of the art performance in the *grandtest* protocol. The model is trained using Binary Cross-Entropy (*BCE*) loss only.

All the baseline methods described are reproducible, and the details about the parameters can be found in our open-source package [4].

### F. Experiments and Results in WMCA dataset

We have tested the baselines and the proposed approach in three different protocols in *WMCA*. The proposed approach is denoted as *MCCNN(BCE+OCCL)-GMM*.

- *MCCNN(BCE+OCCL)-GMM*: Here, the *bonafide* embeddings from the *MCCNN* trained using both the losses are used to train a GMM, and in the evaluation stage, the score from the one class GMM is used as the PAD score.

The results in each protocol are described below.

*1) Experiments in grandtest protocol:* The *grandtest* protocol emulates the known attack scenario. Table II tabulates the results in the *grandtest* protocol. The proposed approach outperforms the feature-based methods by a large margin as expected. The model *MC-RDWT-Haralick-GMM* trained using a one-class model achieves the worse results. It is interesting to note that the *MC-RDWT-Haralick-SVM* model, trained using the same feature as a binary classifier performed much better. This shows one weakness of one-class classifiers in a known attack scenario, as they do not use the known attacks in training. The *MCCNN(BCE)* achieves much better performance as compared to *MC-ResNetPAD*. The *MCCNN(BCE)* trained as a binary classifier achieves the best performance in this protocol. The proposed *MCCNN(BCE+OCCL)-GMM* approach achieves comparable performance to *MCCNN(BCE)*. This indicates that the one class GMM classifier performs on par with the binary classification, provided they are trained with compact feature representations.

*2) Experiments in unseen-2D and unseen-3D protocol:* The *unseen-2D* and *unseen-3D* protocols emulates the unseen attack scenario. The *unseen-3D* is the most challenging protocol since it is trained only on 2D - print and replay attacks and encounters a wide variety of 3D attacks such as silicone masks, fake heads, mannequins, etc. in the *eval* set.

Most of the approaches perform well in the *unseen-2D* protocol. This result is intuitive as these models are trained on challenging 3D attacks, detection of 2D attacks is much easier. Moreover, the 2D attacks can be easily identified in depth, thermal, and infrared channels. Even some feature-based methods perform well in this protocol, with *MC-RDWT-Haralick-GMM* method achieving the best performance. This

shows the advantage of one class model in an unseen attack scenario. The proposed approach *MCCNN(BCE+OCCL)-GMM* and *MCCNN(BCE)* baseline perform comparably in this protocol. Notably, the DeepPixBiS model achieves much worse results in this protocol. This could be because discriminating between *bonafide* and 2D attacks are harder when only RGB information is used.

The *unseen-3D* protocol shows important results. All the baselines show inferior performance when encountered with unseen 3D samples. This shows the failure of binary classifiers in generalizing to challenging unseen attacks. The *MCCNN(BCE)* approach, while being architecturally similar, fails to generalize when trained in the binary classification setting. With the proposed approach, performance improves to 9.7% when the one class GMM is used on the *bonafide* representations. Since the network learns to map the *bonafide* samples to a compact cluster in the feature space, even in the presence of unseen attacks, the decision boundary learned for the *bonafide* class is robust. The unseen attacks map far from the *bonafide* cluster and hence becomes easy to discriminate from *bonafide* samples. This result is encouraging since the network was shown only 2D attacks in training, and still, it manages to achieve good performance against challenging 3D attacks. The ROCs for all the protocols are shown in Fig. 7.

The t-SNE [51] plots of the embeddings for all protocols are shown in Fig. 8. Five frames from each video in the evaluation sets of the protocols are used for this visualization. While the difference between *bonafide* and attacks are clear in the *grandtest* and *unseen-2D*, difference in *unseen-3D* protocol is very evident. It can be clearly seen that the *bonafide* class clusters together and is far from the *bonafide* representation in the embedding space in the *unseen-3D* protocol when the proposed loss is used. Unseen attacks overlaps with *bonafide* embeddings when only *BCE* is used. This clearly demonstrates the effectiveness of the proposed approach for unseen attack detection. The unseen attacks which are overlapping with the *bonafide* region are shown in Fig. 9. It can be seen that some video replay samples and flexible silicone 3D masks get misclassified in unseen-2D and unseen-3D protocols respectively.

*3) Ablation study with channels:* To evaluate the performance of the proposed framework on different set of channels, we perform an ablation study by including a different set of channels. We used only the best performing *MCCNN(BCE+OCCL)-GMM* approach in this ablation study. In all combinations, the gray-scale channel is present since it is used as a reference. This is required as the embedding from the gray-scale part can be used for face recognition as well.

The acronyms for different channels are shown below:

- G: Gray-scale image
- D: Depth image
- I: Infrared channel
- T: Thermal channel

Various combinations of these channels are experimented with, and the results are tabulated in Table IV. It is to be noted that the channels G, D and I come from the same device and T is coming from a different device. Usually, thermal cameras are expensive, compared to RGB-D cameras, and

---

[3]Available from: https://github.com/AlexanderParkin/ChaLearn_liveness_challenge

[4]Source code: https://gitlab.idiap.ch/bob/bob.paper.oneclass_mccnn_2019
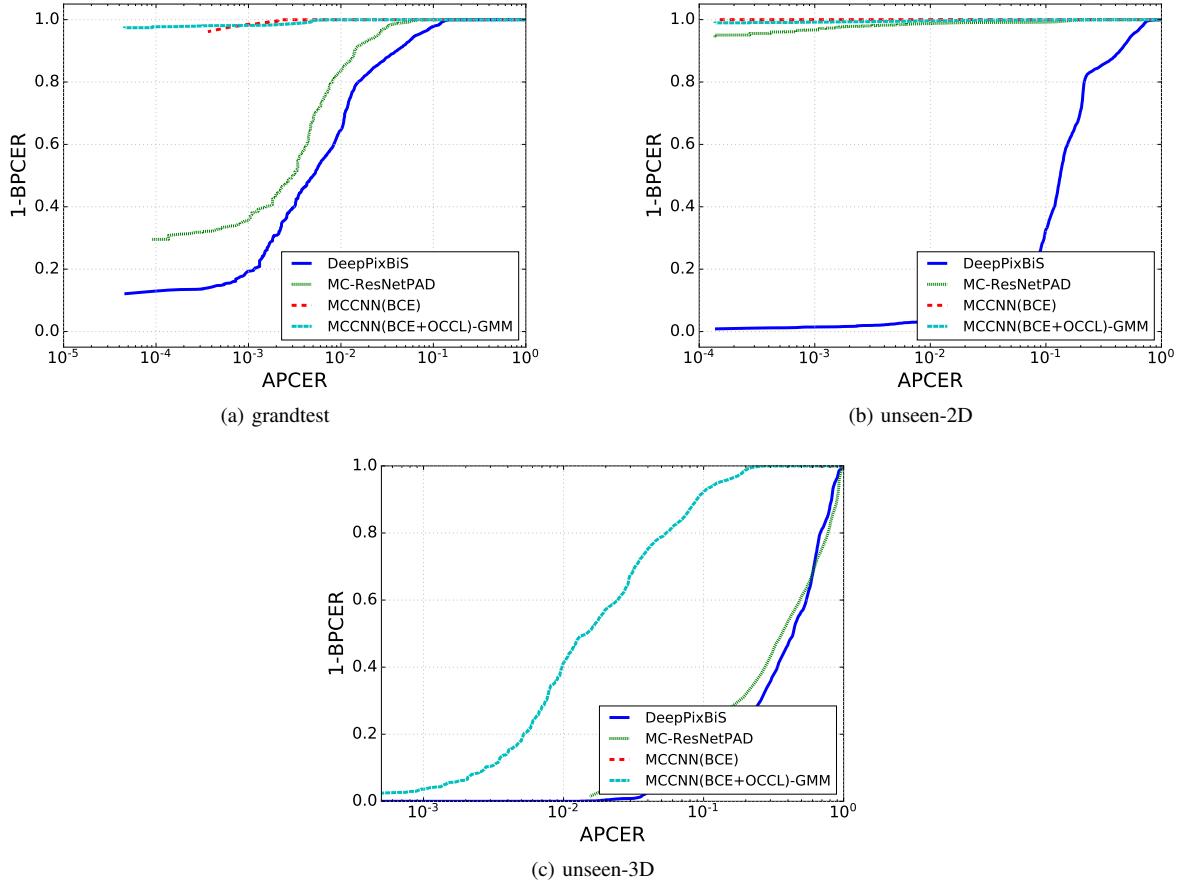
(a) grandtest

(b) unseen-2D

(c) unseen-3D

Fig. 7. DET curves for the *eval* sets of different protocols of *WMCA* dataset a) *grandtest*, b) *unseen-2D*, c) *unseen-3D* protocol.

TABLE II
PERFORMANCE OF THE BASELINE SYSTEMS AND THE PROPOSED METHOD IN *grandtest* PROTOCOL OF *WMCA* DATASET. THE VALUES REPORTED ARE OBTAINED WITH A THRESHOLD COMPUTED FOR BPCER 1% IN *dev* SET.

| Method | dev (%) | | test (%) | | |
|---|---|---|---|---|---|
| | APCER | ACER | APCER | BPCER | ACER |
| MC-RDWT-Haralick-SVM | 3.6 | 2.3 | 5.4 | 1.2 | 3.3 |
| MC-LBP-SVM | 3.6 | 2.3 | 8.5 | 0.6 | 4.6 |
| MC-RDWT-Haralick-GMM | 43.4 | 22.2 | 47.7 | 1.7 | 24.7 |
| DeepPixBiS (RGB only)[10] | 1.0 | 1.0 | 8.2 | 3.7 | 6 |
| MC-ResNetPAD [38] | 3.8 | 2.4 | 3.5 | 1.6 | 2.6 |
| MCCNN(BCE)[14] | 0.4 | 0.7 | 0.5 | 0 | **0.2** |
| **MCCNN(BCE+OCCL)-GMM** | 0.1 | 0.6 | 0.6 | 0.1 | 0.4 |

TABLE III
PERFORMANCE OF THE BASELINE SYSTEMS AND THE PROPOSED METHOD IN **UNSEEN** PROTOCOLS OF *WMCA* DATASET. THE VALUES REPORTED ARE OBTAINED WITH A THRESHOLD COMPUTED FOR BPCER 1% IN *dev* SET.

| Method | unseen-2D | | | unseen-3D | | |
|---|---|---|---|---|---|---|
| | APCER | BPCER | ACER | APCER | BPCER | ACER |
| MC-RDWT-Haralick-SVM | 0.3 | 0.1 | 0.2 | 66.0 | 0.1 | 33.1 |
| MC-LBP-SVM | 40.7 | 0.1 | 20.4 | 38.9 | 0.2 | 19.5 |
| MC-RDWT-Haralick-GMM | 0.0 | 0.2 | **0.1** | 70.8 | 1.9 | 36.4 |
| DeepPixBiS (RGB only)[10] | 77.7 | 0.3 | 39 | 74.7 | 16.3 | 45.5 |
| MC-ResNetPAD [38] | 4.1 | 0.9 | 2.5 | 92.2 | 6.4 | 49.3 |
| MCCNN(BCE)[14] | 0.0 | 1.0 | 0.5 | 62.0 | 0.0 | 31.0 |
| **MCCNN(BCE+OCCL)-GMM** | 0.3 | 0.6 | 0.5 | 15.4 | 3.9 | **9.7** |

(a) grandtest      (b) unseen-2D      (c) unseen-3D
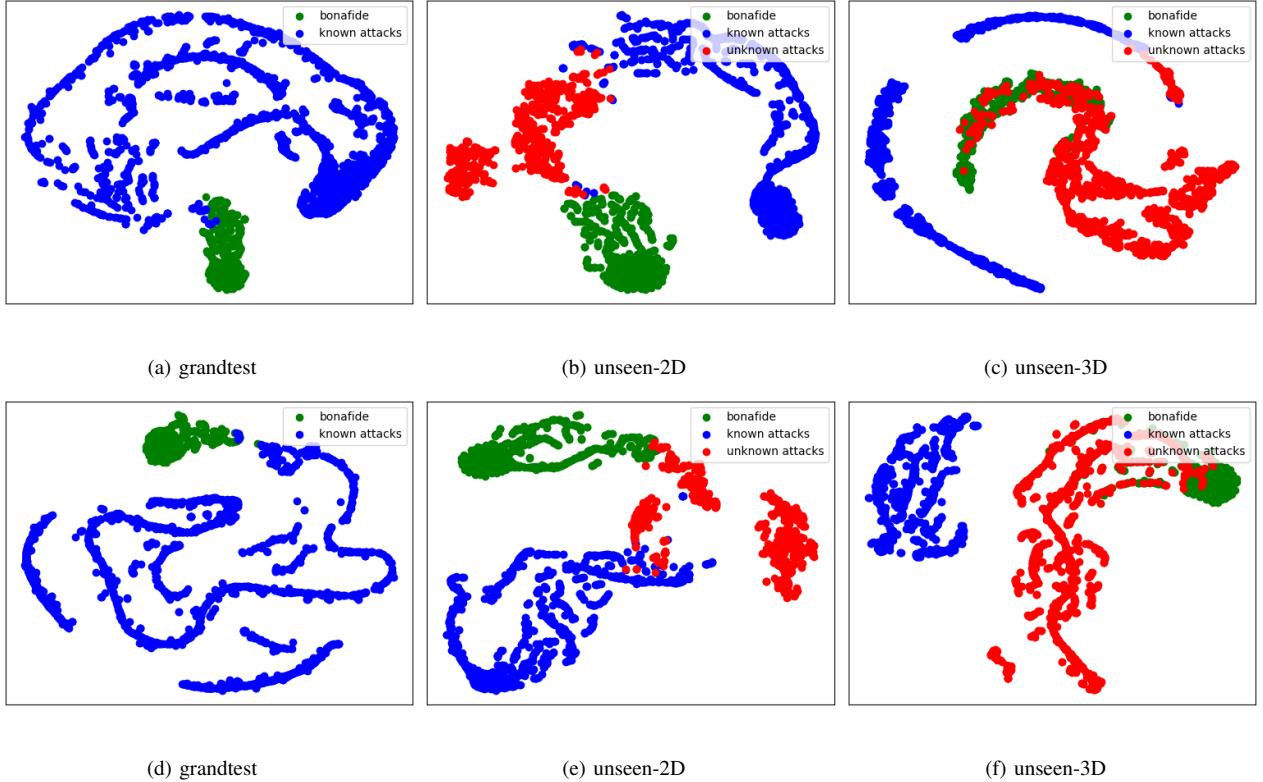
(d) grandtest      (e) unseen-2D      (f) unseen-3D

Fig. 8. t-SNE plots of embeddings in the protocols in *WMCA* dataset. First row (a,b,c) shows the embeddings when only *BCE* loss was used. Second row (d,e,f) shows the embeddings when both the losses are used. Embeddings of both known and unseen attacks are shown in the figures for each protocol. Grandtest protocol contains only known attacks in the test set.
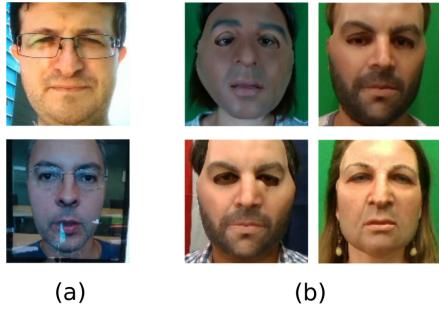


(a)        (b)

Fig. 9. The attack samples which are closer to *bonafide* cluster in a) unseen-2D (Fig.8(E)) and b) unseen-3D ((Fig.8(F))) protocol for the proposed framework.

hence the combinations involving subsets of G, D and I are more interesting from a deployment point of view.

From Table IV, it can be seen that the performance degrades as channels are removed. However, the combination GI achieves reasonable performance while considering the performance-cost ratio. The ROCs for different protocols are shown in Fig. 10.

### G. Experiments and Results in MLFP dataset

We have used only the RGB channel for the experiments since the other channels were not captured simultaneously. For the MCCNN framework and other baselines, 'R', 'G', and 'B' are considered as the different channels in these experiments.

TABLE IV
PERFORMANCE OF THE PROPOSED FRAMEWORK WITH DIFFERENT
COMBINATIONS OF CHANNELS IN ALL PROTOCOLS OF *WMCA* DATASET.
THE VALUES REPORTED ARE OBTAINED WITH A THRESHOLD COMPUTED
FOR BPCER 1% IN *dev* SET.

| Channels | grandtest | unseen-2D | unseen-3D |
|---|---|---|---|
| | ACER | ACER | ACER |
| GDIT | **0.4** | **0.5** | **9.7** |
| GDI | 1.1 | 11.2 | 23.1 |
| GT | 2.2 | 3.2 | 21.5 |
| GD | 2.3 | 49.4 | 45.4 |
| GI | 1.1 | 2.2 | 22.6 |

We have performed the experiments in the three newly created protocols and the results are tabulated in Table V.

From the results in Table V, it can be seen that the CNN based approach outperforms the feature-based approaches. The MCCNN framework, with the addition of the newly proposed loss outperforms the architecture trained with BCE only, showing the effectiveness of the proposed approach.

Even though the proposed approach performs better than the baselines, it is to be noted that the key point of the proposed approach, leveraging multi-channel information, is not utilized here. The architecture is not optimized for PAD in RGB and this experiment is performed only to show the change in performance with the new loss function. Nevertheless, the proposed approach achieves better performance as compared
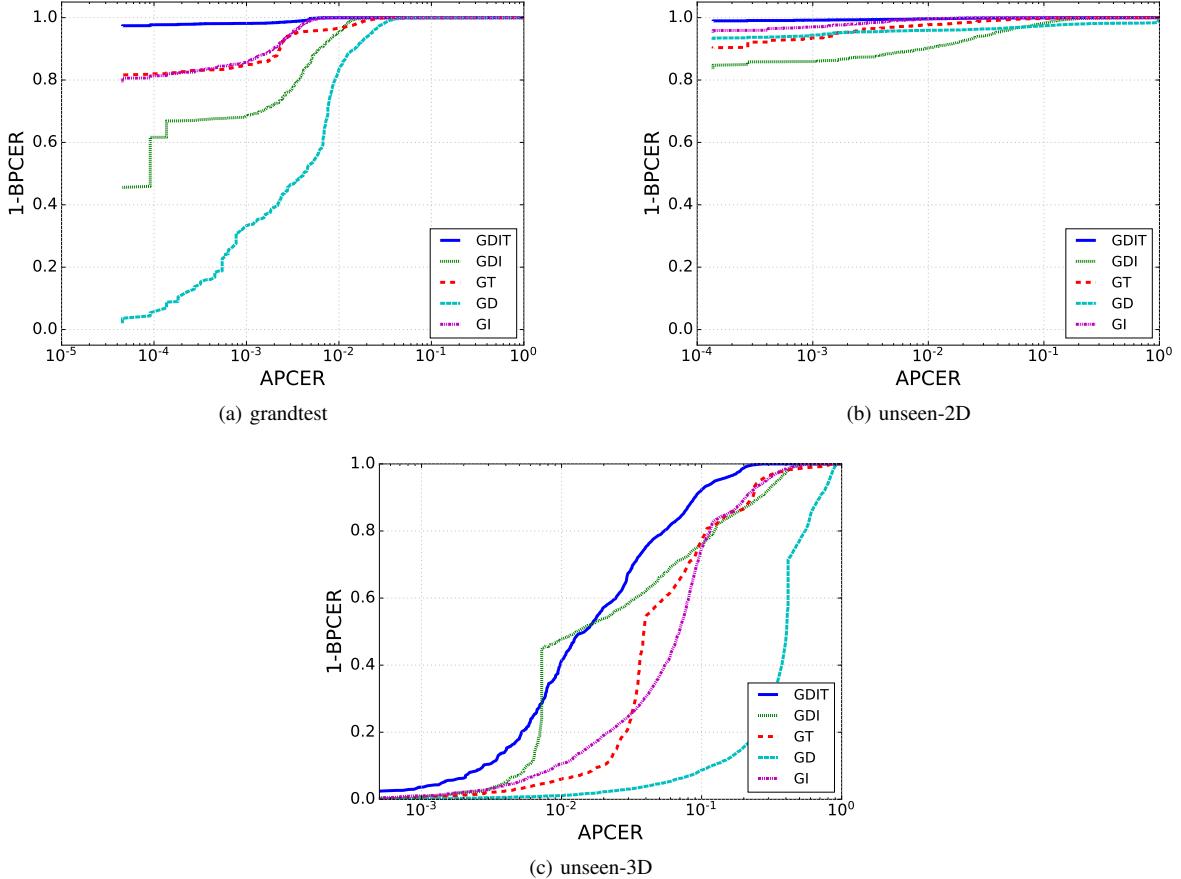
(a) grandtest



(b) unseen-2D



(c) unseen-3D

Fig. 10. Ablation study with different combination of channels, DET curves for the *eval* sets of different protocols of *WMCA* dataset a) *grandtest*, b) *unseen-2D*, c) *unseen-3D* protocol.

to the baselines in all the protocols.

TABLE V
PERFORMANCE OF THE PROPOSED FRAMEWORK IN THE PROTOCOLS IN *MLFP* DATASET. ONLY RGB CHANNEL WAS USED IN THIS EXPERIMENTS. THE VALUES REPORTED ARE THE EER IN THE *evaluation* SET.

| Algorithm | grandtest | unseen print | unseen mask |
|---|---|---|---|
| MC-RDWT-Haralick-SVM | 9.8 | 12.0 | 32.2 |
| MC-LBP-SVM | 6.3 | 27.1 | 9.3 |
| MC-RDWT-Haralick-GMM | 27.4 | 40.8 | 21.5 |
| DeepPixBiS (RGB only)[10] | 6.3 | 24.8 | 17.5 |
| MCCNN (BCE) | 5.5 | 9.2 | 5.2 |
| **MCCNN (BCE+OCCL)-GMM** | **1.2** | **3.3** | **3.4** |

### H. Experiments and Results in SiW-M dataset

Table VI shows the performance of the proposed framework, again in RGB only scenario. CNN based methods performs much better than feature based methods in this case. It can be seen that the proposed approach achieves better performance as compared to baseline methods. The performance of the *MCCNN (BCE+OCCL)-GMM* is better compared to the *MCCNN(BCE)* model. It can be seen that the addition of the new loss function makes the classification of unseen attacks more accurate.

### I. Cross-database evaluations

As we could not perform cross-database evaluation between a multi-channel database and an RGB only database, we used only the RGB channels from two datasets for the cross-database evaluation. We have selected WMCA and SiW-M datasets as they are relatively large and consist of a wide variety of attacks.

From Table VII, it can be seen that the MCCNN model with and without the new loss performs comparably. In general, the performance in the cross-database setting is poor for all the models. The poor performance could be due to the disparity in acquisition conditions and the attack types. A wider variety of attacks makes it difficult for the classifier to identify attacks using RGB channels alone. The cross-database performance with this wide variety of attacks seems more challenging as compared to typical cross-database evaluations using only 2D attacks. Using multiple channels [14] may alleviate these issues. This also points to the limitation of RGB only methods while dealing with a wide variety of attacks.

### V. DISCUSSIONS

From the experiments in *WMCA* database, it can be clearly seen that CNN based method outperforms the feature-based methods by a large margin. While comparing the *MC-CNN(BCE)* method to the proposed method, the performance

TABLE VI
PERFORMANCE OF THE PROPOSED FRAMEWORK IN THE LEAVE ONE OUT PROTOCOLS IN *SiW-M* DATASET. ONLY RGB CHANNEL WAS PRESENT IN THIS DATASET.

| Methods | Metrics (%) | Replay | Print | Mask Attacks | | | | | Makeup Attacks | | | Partial Attacks | | | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Half | Silicone | Trans. | Paper | Manne. | Obfusc. | Imperson. | Cosmetic | Funny Eye | Paper Glasses | Partial Paper | |
| MC-RDWT-Haralick-SVM | APCER | 19.80 | 19.15 | 30.76 | 28.15 | 33.35 | 0.29 | 4.50 | 68.91 | 0.00 | 35.20 | 53.12 | 34.53 | 3.49 | 25.4 ± 20.8 |
| | BPCER | 14.50 | 13.89 | 14.66 | 16.83 | 15.38 | 16.68 | 15.88 | 16.03 | 16.53 | 16.37 | 14.58 | 14.47 | 15.73 | 15.5 ± 0.9 |
| | ACER | 17.15 | 16.52 | 22.71 | 22.49 | 24.37 | 8.49 | 10.19 | 42.47 | 8.26 | 25.79 | 33.85 | 24.50 | 9.61 | 20.4 ± 10.3 |
| | EER | 16.88 | 16.53 | 21.80 | 20.73 | 21.94 | 7.34 | 9.88 | 32.56 | 2.37 | 23.51 | 31.72 | 21.94 | 10.05 | 18.2 ± 9.0 |
| MC-LBP-SVM | APCER | 10.77 | 12.91 | 10.28 | 35.19 | 37.78 | 0.59 | 6.50 | 96.09 | 0.00 | 26.00 | 40.91 | 35.51 | 2.73 | 24.2 ± 26.3 |
| | BPCER | 22.90 | 22.18 | 22.48 | 22.33 | 23.13 | 23.70 | 23.59 | 22.79 | 23.93 | 22.90 | 19.92 | 21.11 | 23.74 | 22.6 ± 1.1 |
| | ACER | 16.83 | 17.54 | 16.38 | 28.76 | 30.46 | 12.15 | 15.04 | 59.44 | 11.97 | 24.45 | 30.41 | 28.31 | 13.24 | 23.4 ± 12.9 |
| | EER | 15.96 | 16.83 | 16.87 | 28.51 | 29.77 | 10.54 | 12.75 | 52.60 | 1.90 | 24.61 | 28.32 | 26.76 | 11.29 | 21.2 ± 12.6 |
| Auxiliary [19] | APCER | 23.7 | 7.3 | 27.7 | 18.2 | 97.8 | 8.3 | 16.2 | 100.0 | 18.0 | 16.3 | 91.8 | 72.2 | 0.4 | 38.3 ± 37.4 |
| | BPCER | 10.1 | 6.5 | 10.9 | 11.6 | 6.2 | 7.8 | 9.3 | 11.6 | 9.3 | 7.1 | 6.2 | 8.8 | 10.3 | 8.9 ± 2.0 |
| | ACER | 16.8 | 6.9 | 19.3 | 14.9 | 52.1 | 8.0 | 12.8 | 55.8 | 13.7 | 11.7 | 49.0 | 40.5 | 5.3 | 23.6 ± 18.5 |
| | EER | 14.0 | 4.3 | 11.6 | 12.4 | 24.6 | 7.8 | 10.0 | 72.3 | 10.1 | 9.4 | 21.4 | 18.6 | 4.0 | 17.0 ± 17.7 |
| Deep Tree Network [28] | APCER | 1.0 | 0.0 | 0.7 | 24.5 | 58.6 | 0.5 | 3.8 | 73.2 | 13.2 | 12.4 | 17.0 | 17.0 | 0.2 | 17.1 ± 23.3 |
| | BPCER | 18.6 | 11.9 | 29.3 | 12.8 | 13.4 | 8.5 | 23.0 | 11.5 | 9.6 | 16.0 | 21.5 | 22.6 | 16.8 | 16.6 ± 6.2 |
| | ACER | 9.8 | 6.0 | 15.0 | 18.7 | 36.0 | 4.5 | 7.7 | 48.1 | 11.4 | 14.2 | 19.3 | 19.8 | 8.5 | 16.8 ± 11.1 |
| | EER | 10.0 | 2.1 | 14.4 | 18.6 | 26.5 | 5.7 | 9.6 | 50.2 | 10.1 | 13.2 | 19.8 | 20.5 | 8.8 | 16.1 ± 12.2 |
| DeepPixBiS [10] | APCER | 19.18 | 8.97 | 1.74 | 21.30 | 60.68 | 0.00 | 1.00 | 100.00 | 0.00 | 26.90 | 64.66 | 77.52 | 0.29 | 29.4± 34.4 |
| | BPCER | 8.70 | 7.63 | 11.03 | 11.76 | 10.27 | 8.85 | 8.63 | 10.53 | 11.60 | 10.99 | 10.31 | 10.23 | 7.10 | 9.8± 1.4 |
| | ACER | 13.94 | 8.30 | 6.38 | 16.53 | 35.47 | 4.43 | 4.81 | 55.27 | 5.80 | 18.95 | 37.48 | 43.87 | 3.69 | 19.6± 17.4 |
| | EER | 11.68 | 7.94 | 7.22 | 15.04 | 21.30 | 3.78 | 4.52 | 26.49 | 1.23 | 14.89 | 23.28 | 18.90 | 4.82 | 12.3± 8.2 |
| MCCNN (BCE) | APCER | 38.93 | 30.60 | 7.85 | 20.00 | 32.56 | 0.00 | 2.00 | 70.65 | 0.00 | 29.00 | 46.69 | 57.32 | 23.20 | 27.6 ± 22.1 |
| | BPCER | 7.10 | 6.45 | 7.48 | 10.04 | 12.56 | 8.59 | 10.04 | 9.96 | 11.72 | 11.37 | 12.75 | 7.71 | 9.89 | 9.6 ± 2.2 |
| | ACER | 23.01 | 18.52 | 7.66 | 15.02 | 22.56 | 4.29 | 6.02 | 40.31 | 5.86 | 20.19 | 29.72 | 32.52 | 16.54 | 18.6 ± 11.1 |
| | EER | 17.08 | 11.83 | 7.56 | 12.82 | 16.09 | 0.71 | 6.85 | 25.94 | 2.29 | 16.30 | 18.90 | 22.82 | 13.13 | 13.2 ± 7.4 |
| **MCCNN (BCE+OCCL)-GMM** | APCER | 11.79 | 9.53 | 3.12 | 3.70 | 39.20 | 0.00 | 3.12 | 44.57 | 0.00 | 21.60 | 19.34 | 35.55 | 0.00 | 14.7 ± 15.9 |
| | BPCER | 13.44 | 16.15 | 16.26 | 20.23 | 11.11 | 13.74 | 8.66 | 15.23 | 12.67 | 10.42 | 14.31 | 18.40 | 27.33 | 15.2 ± 4.8 |
| | ACER | 12.61 | 12.84 | 9.69 | 11.97 | 25.16 | 6.87 | 5.89 | 29.90 | 6.34 | 16.01 | 16.83 | 26.97 | 13.66 | 14.9 ± 7.8 |
| | EER | 12.82 | 12.94 | 11.33 | 13.70 | 13.47 | 0.56 | 5.60 | 22.17 | 0.59 | 15.14 | 14.40 | 23.93 | 9.82 | **12.0 ± 6.9** |

TABLE VII
THE RESULTS FROM THE CROSS-DATABASE TESTING BETWEEN WMCA AND SiW-M DATASETS USING THE GRANDTEST PROTOCOL, ONLY RGB CHANNELS WERE USED IN THIS EXPERIMENT.

| Method | trained on WMCA | | trained on SiW-M | |
|---|---|---|---|---|
| | tested on WMCA | tested on SiW-M | tested on SiW-M | tested on WMCA |
| MC-RDWT-Haralick-SVM | 14.6 | 29.6 | 15.1 | 45.3 |
| MC-LBP-SVM | 26.6 | 45.5 | 19.6 | 38.6 |
| MC-RDWT-Haralick-GMM | 27.9 | 34.0 | 25.5 | 43.6 |
| DeepPixBiS | 7.5 | 49.1 | 14.7 | 44.4 |
| MCCNN (BCE) | 12.1 | 34.0 | 9.9 | 42.3 |
| **MCCNN (BCE+OCCL)-GMM** | 12.3 | 31.9 | 9.5 | 41.8 |

is comparable in the known attack scenario. This indicates that the proposed One-Class GMM based approach performs par with binary classification, thanks to the embedding learned with the proposed loss function. Most of the approaches perform well in the *unseen-2D* protocol since it can be clearly discriminated in many channels. Moreover, it shows that if the network is trained in challenging attacks, simpler attacks are easy to detect. While the performance is comparable in *grandtest* and *unseen-2D* protocols, the proposed method achieves a large performance boost in the most challenging *unseen-3D* protocol. The proposed loss function forces the network to learn a compact representation for *bonafide* samples in the feature space. Both known and unknown attacks get mapped far from the *bonafide* cluster in the feature space. The decision boundary learned by the one class model seems to be robust in identifying both seen and unseen attacks in such a scenario. This result is significant for several reasons. It is to be noted that in the *unseen-3D* protocol, the network is trained with only 2D attacks, i.e., prints and replays. The proposed method achieves excellent performance in a test set consisting of challenging 3D attacks such as custom silicone masks, paper masks, mannequins, etc. The real-world implications of this approach is very promising. The proposed method can be used to develop robust PAD systems without the requirement of having to manufacture costly presentation attacks. The models can be trained on easy to obtain attacks based on availability. The proposed framework utilizes the available (known) attack categories to learn a robust representation to facilitate known and unseen attack detection. It is to be noted that the compact representation is made possible by the joint multi-channel representation used.

In practical deployment scenarios, it may not be possible to use all the four channels due to the computational or cost constraints. In such a situation, models trained on available channels can be selected based on the performance-cost ratio by sub selecting the channels. The results from the ablation study presented in Table IV can be used to determine which channels should be used in such cases.

Similarly, the experiments in *MLFP* and *SiW-M* database also shows CNN based methods outperform feature-based baselines. Although we have not used multi-channel information in the experiments, the experiment results showcase the performance improvement with the new loss function. Using the proposed framework together with network backbones designed specifically for RGB PAD might improve the results.

The cross-database performance shows the limitations of the RGB channel when tested with a wide variety of attacks. The performance of the baselines as well as the proposed approach is poor when only RGB data is used. This shows the challenging nature of RGB only PAD while considering a wide variety of attacks. The usage of multiple channels as done in *WMCA* dataset might improve the performance.

## VI. CONCLUSIONS

Face presentation attack detection is often considered as a binary classification task which results in over-fitting to the known attacks leading to poor generalization against unseen attacks. In this paper, we address this problem with a new

framework using a one-class classifier. A novel loss function is formulated, which forces the CNN to learn a compact yet discriminative representation for the face images. The *bonafide* samples form a compact cluster in the feature space, thanks to the proposed loss function. A decision boundary around the *bonafide* representation can be obtained using a one-class model. Both known and unknown attacks map far from the *bonafide* cluster in the feature space which can be classified by the one-class model. The proposed framework introduces a new way to learn a robust PAD system from *bonafide* and available (known) attack classes. The proposed system has been evaluated in the challenging *WMCA*, *MLFP*, and *SiW-M* databases and was shown to outperform the baseline feature-based and CNN based methods in both known and unseen attack scenarios. The drastic improvement in the performance in *unseen-3D* protocol in *WMCA* shows the robustness of the proposed approach against unseen attacks, thanks to the multi-channel information. The proposed method also shows improvement even when used together with RGB channels alone. The source code and protocols to reproduce the results are made available publicly to enable further extensions of the proposed framework.

## REFERENCES

[1] A. K. Jain and S. Z. Li, *Handbook of face recognition*. Springer, 2011.

[2] E. Learned-Miller, G. B. Huang, A. RoyChowdhury, H. Li, and G. Hua, "Labeled faces in the wild: A survey," in *Advances in face detection and facial image analysis*. Springer, 2016, pp. 189–248.

[3] S. Marcel, M. S. Nixon, J. Fierrez, and N. Evans, *Handbook of biometric anti-spoofing : Presentation attack detection*. Editors: Marcel, S., Nixon, M.S., Fierrez, J., Evans, N. (Eds.); Springer International Publishing, 2018, 2nd ed.; ISBN: 978-3319926261, 09 2018. [Online]. Available: http://www.eurecom.fr/publication/5667

[4] ISO/IEC JTC 1/SC 37 Biometrics, "Information technology International Organization for Standardization," International Organization for Standardization, ISO Standard, Feb. 2016.

[5] Z. Boulkenafet, J. Komulainen, and A. Hadid, "Face anti-spoofing based on color texture analysis," in *Image Processing (ICIP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 2636–2640.

[6] J. Määttä, A. Hadid, and M. Pietikäinen, "Face spoofing detection from single images using micro-texture analysis," in *Biometrics (IJCB), 2011 international joint conference on*. IEEE, 2011, pp. 1–7.

[7] A. Anjos and S. Marcel, "Counter-measures to photo attacks in face recognition: a public database and a baseline," in *Biometrics (IJCB), 2011 international joint conference on*. IEEE, 2011, pp. 1–7.

[8] R. Ramachandra and C. Busch, "Presentation attack detection methods for face recognition systems: a comprehensive survey," *ACM Computing Surveys (CSUR)*, vol. 50, no. 1, p. 8, 2017.

[9] G. Heusch and S. Marcel, "Remote blood pulse analysis for face presentation attack detection," in *Handbook of Biometric Anti-Spoofing*. Springer, 2019, pp. 267–289.

[10] A. George and S. Marcel, "Deep pixel-wise binary supervision for face presentation attack detection," *International Conference on Biometrics*, 2019.

[11] R. Raghavendra, K. B. Raja, S. Venkatesh, and C. Busch, "Extended multispectral face presentation attack detection: An approach based on fusing information from individual spectral bands," in *Information Fusion (Fusion), 2017 20th International Conference on*. IEEE, 2017, pp. 1–6.

[12] H. Steiner, A. Kolb, and N. Jung, "Reliable face anti-spoofing using multispectral swir imaging," in *Biometrics (ICB), 2016 International Conference on*. IEEE, 2016, pp. 1–8.

[13] S. Bhattacharjee and S. Marcel, "What you can't see can help you–extended-range imaging for 3d-mask presentation attack detection," in *Proceedings of the 16th International Conference on Biometrics Special Interest Group.*, no. EPFL-CONF-231840. Gesellschaft fuer Informatik eV (GI), 2017.

[14] 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*, pp. 1–1, 2019.

[15] A. George and S. Marcel, "Can your face detector do anti-spoofing? face presentation attack detection with a multi-channel face detector," *Idiap Research Report, Idiap-RR-12-2020*, 2020.

[16] G. Heusch, A. George, D. Geissbühler, Z. Mostaani, and S. Marcel, "Deep models and shortwave infrared information to detect face presentation attacks," *IEEE Transactions on Biometrics, Behavior, and Identity Science (T-BIOM)*, 2020.

[17] O. Nikisins, A. George, and S. Marcel, "Domain adaptation in multi-channel autoencoder based features for robust face anti-spoofing," in *2019 International Conference on Biometrics (ICB)*. IEEE, 2019, pp. 1–8.

[18] S. Bhattacharjee, A. Mohammadi, and S. Marcel, "Spoofing deep face recognition with custom silicone masks," *Biometrics Theory, Applications and Systems (BTAS), 2018 IEEE 9th International Conference on*, 2018.

[19] 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*, 2018, pp. 389–398.

[20] Y. Atoum, Y. Liu, A. Jourabloo, and X. Liu, "Face anti-spoofing using patch and depth-based cnns," in *Biometrics (IJCB), 2017 IEEE International Joint Conference on*. IEEE, 2017, pp. 319–328.

[21] R. Shao, X. Lan, and P. C. Yuen, "Deep convolutional dynamic texture learning with adaptive channel-discriminability for 3d mask face anti-spoofing," in *Biometrics (IJCB), 2017 IEEE International Joint Conference on*. IEEE, 2017, pp. 748–755.

[22] S. R. Arashloo, J. Kittler, and W. Christmas, "An anomaly detection approach to face spoofing detection: A new formulation and evaluation protocol," *IEEE Access*, vol. 5, pp. 13 868–13 882, 2017.

[23] O. Nikisins, A. Mohammadi, A. Anjos, and S. Marcel, "On effectiveness of anomaly detection approaches against unseen presentation attacks in face anti-spoofing," in *The 11th IAPR International Conference on Biometrics (ICB 2018)*, no. EPFL-CONF-233583, 2018.

[24] F. Xiong and W. AbdAlmageed, "Unknown presentation attack detection with face rgb images," in *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2018, pp. 1–9.

[25] P. Perera and V. M. Patel, "Learning deep features for one-class classification," *IEEE Transactions on Image Processing*, vol. 28, no. 11, pp. 5450–5463, 2019.

[26] S. Fatemifar, M. Awais, S. R. Arashloo, and J. Kittler, "Combining multiple one-class classifiers for anomaly based face spoofing attack detection," in *International Conference on Biometrics (ICB)*, 2019.

[27] D. Pérez-Cabo, D. Jiménez-Cabello, A. Costa-Pazo, and R. J. López-Sastre, "Deep anomaly detection for generalized face anti-spoofing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

[28] Y. Liu, J. Stehouwer, A. Jourabloo, and X. Liu, "Deep tree learning for zero-shot face anti-spoofing," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[29] A. Jaiswal, S. Xia, I. Masi, and W. AbdAlmageed, "Ropad: Robust presentation attack detection through unsupervised adversarial invariance," *arXiv preprint arXiv:1903.03691*, 2019.

[30] S. Mehta, A. Uberoi, A. Agarwal, M. Vatsa, and R. Singh, "Crafting a panoptic face presentation attack detector."

[31] J. J. Engelsma and A. K. Jain, "Generalizing fingerprint spoof detector: Learning a one-class classifier," *arXiv preprint arXiv:1901.03918*, 2019.

[32] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[33] N. Erdogmus and S. Marcel, "Spoofing face recognition with 3d masks," *IEEE transactions on information forensics and security*, vol. 9, no. 7, pp. 1084–1097, 2014.

[34] T. I. Dhamecha, A. Nigam, R. Singh, and M. Vatsa, "Disguise detection and face recognition in visible and thermal spectrums," in *Biometrics (ICB), 2013 International Conference on*. IEEE, 2013, pp. 1–8.

[35] A. Agarwal, D. Yadav, N. Kohli, R. Singh, M. Vatsa, and A. Noore, "Face presentation attack with latex masks in multispectral videos," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, July 2017, pp. 275–283.

[36] G. Wang, C. Lan, H. Han, S. Shan, and X. Chen, "Multi-modal face presentation attack detection via spatial and channel attentions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

[37] S. Zhang, X. Wang, A. Liu, C. Zhao, J. Wan, S. Escalera, H. Shi, Z. Wang, and S. Z. Li, "Casia-surf: A dataset and benchmark for large-scale multi-modal face anti-spoofing," *arXiv preprint arXiv:1812.00408*, 2018.

[38] A. Parkin and O. Grinchuk, "Recognizing multi-modal face spoofing with face recognition networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

[39] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *European conference on computer vision*. Springer, 2016, pp. 499–515.

[40] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2. IEEE, 2006, pp. 1735–1742.

[41] C. Qi and F. Su, "Contrastive-center loss for deep neural networks," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 2851–2855.

[42] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.

[43] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 532–539.

[44] C. Leys, C. Ley, O. Klein, P. Bernard, and L. Licata, "Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median," *Journal of Experimental Social Psychology*, vol. 49, no. 4, pp. 764–766, 2013.

[45] X. Wu, R. He, Z. Sun, and T. Tan, "A light cnn for deep face representation with noisy labels," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 11, pp. 2884–2896, 2018.

[46] T. de Freitas Pereira, A. Anjos, and S. Marcel, "Heterogeneous face recognition using domain specific units," *IEEE Transactions on Information Forensics and Security*, 2018.

[47] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B*, vol. 39, no. 1, pp. 1–38, 1977.

[48] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[49] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.

[50] R. M. Haralick, "Statistical and structural approaches to texture," *Proceedings of the IEEE*, vol. 67, no. 5, pp. 786–804, 1979.

[51] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.

**Anjith George** has received his Ph.D. and M-Tech degree from the Department of Electrical Engineering, Indian Institute of Technology (IIT) Kharagpur, India in 2012 and 2018 respectively. After Ph.D, he worked in Samsung Research Institute as a machine learning researcher. Currently, he is a post-doctoral researcher in the biometric security and privacy group at Idiap Research Institute, focusing on developing face presentation attack detection algorithms. His research interests are real-time signal and image processing, embedded systems, computer vision, machine learning with a special focus on Biometrics.

**Sébastien Marcel** received the Ph.D. degree in signal processing from Université de Rennes I, Rennes, France, in 2000 at CNET, the Research Center of France Telecom (now Orange Labs). He is currently interested in pattern recognition and machine learning with a focus on biometrics security. He is a Senior Researcher at the Idiap Research Institute (CH), where he heads a research team and conducts research on face recognition, speaker recognition, vein recognition, and presentation attack detection (anti-spoofing). He is a Lecturer at the Ecole Polytechnique Fédérale de Lausanne (EPFL) where he teaches a course on "Fundamentals in Statistical Pattern Recognition." He is an Associate Editor of IEEE Signal Processing Letters. He has also served as Associate Editor of IEEE Transactions on Information Forensics and Security, co-editor of the "Handbook of Biometric Anti-Spoofing," Guest Editor of the IEEE Transactions on Information Forensics and Security Special Issue on "Biometric Spoofing and Countermeasures," and co-editor of the IEEE Signal Processing Magazine Special Issue on "Biometric Security and Privacy." He was the Principal Investigator of international research projects including MOBIO (EU FP7 Mobile Biometry), TABULA RASA (EU FP7 Trusted Biometrics under Spoofing Attacks), and BEAT (EU FP7 Biometrics Evaluation and Testing).