

# Adversarially Robust Deepfake Detection via Adversarial Feature Similarity Learning

Sarwar Khan<sup>1,2,3</sup>, Jun-Cheng Chen<sup>1,2</sup>, Wen-Hung Liao<sup>2,3</sup>, and Chu-Song Chen<sup>4</sup>

<sup>1</sup> Research Center for Information Technology Innovation, Academia Sinica

<sup>2</sup> Social Networks Human-Centered Computing, TIGP, Academia Sinica

<sup>3</sup> Computer Science, National Chengchi University

<sup>4</sup> Computer Science and Information Engineering, National Taiwan University  
say2sarwar@gmail.com, pullpull@citi.sinica.edu.tw, whliao@cs.nccu.edu.tw,  
chusong@csie.ntu.edu.tw

**Abstract.** Deepfake technology has raised concerns about the authenticity of digital content, necessitating the development of effective detection methods. However, the widespread availability of deepfakes has given rise to a new challenge in the form of adversarial attacks. Adversaries can manipulate deepfake videos with small, imperceptible perturbations that can deceive the detection models into producing incorrect outputs. To tackle this critical issue, we introduce Adversarial Feature Similarity Learning (AFSL), which integrates three fundamental deep feature learning paradigms. By optimizing the similarity between samples and weight vectors, our approach aims to distinguish between real and fake instances. Additionally, we aim to maximize the similarity between both adversarially perturbed examples and unperturbed examples, regardless of their real or fake nature. Moreover, we introduce a regularization technique that maximizes the dissimilarity between real and fake samples, ensuring a clear separation between these two categories. With extensive experiments on popular deepfake datasets, including FaceForensics++, FaceShifter, and DeeperForensics, the proposed method outperforms other standard adversarial training-based defense methods significantly. This further demonstrates the effectiveness of our approach to protecting deepfake detectors from adversarial attacks.

**Keywords:** Adversarial attack · Adversarial training · Deepfake video detection · Forgery detector.

## 1 Introduction

Deepfakes are synthetic videos in which a person’s face is altered to resemble a different individual, resulting in the production of highly realistic footage depicting events that never actually took place [7]. Deepfake technology has captivated and alarmed society by offering the ability to create remarkably convincing and misleading media, which in turn threatens the authenticity of digital content. In response to these concerns, researchers have diligently worked to develop effective deepfake detection methods [1,8,15,16,37,38,47].

Deepfake detectors have shown promising performance under normal conditions, accurately identifying manipulated videos. However, a new challenge has emerged in the form of adversarial attacks, where small and imperceptible perturbations can deceive the detection models into producing incorrect outputs [10,14,17,18,28,34]. An adversarial example is a manipulated input intentionally designed to deceive a classification model [30]. Adversarial deepfakes [18] leverage the pre-softmax layer to compute the loss and iteratively calculate gradients, allowing for the creation of adversarial fakes that can successfully evade detection. Statistical consistency attack (StatAttack) robust[17] looks into statistical consistency between real and fake and uses degradation techniques to create transferable deepfake adversarial attacks. This poses a significant threat to the reliability and effectiveness of deepfake detection systems, as it undermines their ability to distinguish between genuine and manipulated content.

Adversarial training [30] tackles adversarial attacks through a *min-max* optimization but often at the cost of reduced performance on normal inputs. To address this, TRADES [46] is a surrogate loss for adversarial training utilizing cross-entropy loss supervising and the distance loss between the features of clean and adversarial examples as the regularization. However, training deepfake detectors for robustness remains challenging due to the inclusion of fake images. Our study is motivated by the recognition that adversaries can exploit misclassification, aiming to develop effective strategies for adversarially robust deepfake detection.

In pursuit of this objective, we develop an Adversarial Feature Similarity Learning (AFSL) objective function that optimizes similarity across three fundamental paradigms of deep feature learning. First, we optimize the similarity between samples and weight vectors, specifically focusing on differentiating between real and fake instances. Secondly, our objective is to maximize the similarity between samples, considering both adversarially perturbed examples and unperturbed examples, where the perturbed instances can be either real or fake. Finally, we introduce a regularization approach that aims to maximize the dissimilarity between real and fake samples, ensuring a clear separation and distinct representation of these two categories. This comprehensive approach enables effective deepfake detection by enhancing discrimination between real and fake content and mitigating the impact of adversarial perturbations. We conduct extensive experiments on the FaceForensics++ [37], FaceShifter [25], and DeepForensics [19] datasets, evaluating the performance of our proposed method. Impressively, our method outperforms widely used adversarial training-based defense methods by a significant margin. This demonstrates the effectiveness of our approach to help deepfake detectors fight against various adversarial attacks.

## 2 Related work

### 2.1 Deepfake Creation and Detection

The development of Generative Adversarial Networks (GANs) and their diverse variants has yielded remarkable outcomes in image generation and manipulation,

consequently facilitating the emergence of deepfake technology. By leveraging GANs, deepfake has enabled the creation of fabricated images or videos across various categories. The current deepfake generation techniques include various approaches, such as complete face synthesis [21], face identity swap [11], and face manipulation [12]. The utilization of these generation methods by malicious applications can greatly jeopardize public information security. Nonetheless, it is crucial to recognize that the misuse of deepfake technology raises additional concerns regarding security and privacy, extending to sensitive areas such as politics, religion, and pornography [40,44]. Meanwhile study conducted in Thailand [39] explores Thai perspectives on deepfake AI images, emphasizing both creative interests and potential risks. It suggests the Thai government take a proactive role in regulating and raising awareness to harness the technology’s creative potential while addressing concerns related to data protection and image copyright.

To mitigate the potential misuse of deepfake technologies, various Deep Neural Network (DNN) methods have been proposed for detecting deepfake inputs. Deepfake detection primarily involves the binary classification of distinguishing between fake and real inputs. In the realm of deepfake detection methods, some methods focus on extracting spatial information [8,14,15,16,28,47], whereas others delve into analyzing the differences in frequency information [9] between fake and real inputs. LipForensics [16] employs a lips extraction technique from facial images and leverages a combination of a pretrained feature extractor and a temporal convolutional network to train an effective deepfake detector. FTCN [47] proposed exceptional generalization across various manipulation scenarios by enforcing a uniform spatial convolutional kernel size of one. RealForensics [15] aims to enhance forgery detection performance and improve generalization across different datasets by leveraging real talking faces through self-supervision using spatiotemporal features. These methods achieve remarkable detection results within their respective experimental configurations by harnessing the formidable feature extraction capabilities of DNNs.

## 2.2 Adversarial Examples

Adversarial examples exploit the vulnerability of deep learning models by intentionally designing inputs that cause the models to make mistakes or misclassify data [3]. Gradient-based adversarial attacks are extremely effective against deep learning models in image [3,26,27,30,23], video [32,43], and audio [5,35,36] domain. FGSM [13], PGD [30], CW [3], and StatAttack [17] represent potent adversarial attack techniques that employ distinct optimization strategies to generate perturbations capable of deceiving classification models. Similar to other deep learning models, deepfake detectors are vulnerable to adversarial attacks, making them a significant threat within the realm of deepfakes.

Adversarial deepfakes [18] present a robust white-box attack (RWB) setting, where the attacker possesses full access to the model, as well as a robust black-box (RBB) attack setting, where the attacker lacks access to the model. These settings are achieved by utilizing the pre-softmax layer and employing diverse transformations. Likewise, Neekhara et al. [34] investigate the transferability

of adversarial attacks in forgery detectors and propose a universal attack approach, demonstrating the effectiveness of adversarial examples across different models and architectures. In addition, Statistical Consistency Attack (StatAttack) [17] introduces a transferable approach that leverages adversarial statistical consistency through the minimization of a distribution-aware loss, enabling it to circumvent deepfake detectors effectively.

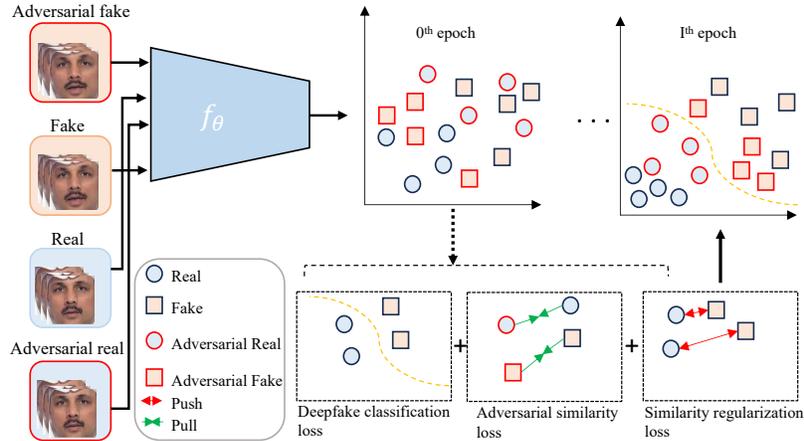
A multitude of defense strategies have been proposed to combat adversarial examples in both image and video domains, adversarial training has exhibited commendable performance in mitigating the impact of such attacks [30,46,29,20,45]. Although adversarial training has shown effectiveness against various forms of adversarial attacks, it often comes at the cost of reduced performance on unperturbed data. To address this, Deep image prior (DIP) [10] utilizes a GAN network to remove perturbations from the input, making the model robust against adversarial attacks, and also incorporates regularization techniques as a defense strategy in deepfake detection. To handle another growing concern of deepfake technology exacerbated by the use of face masks during the pandemic, Alnaim et al. [2] propose a novel deepfake face mask dataset and detection model with identifying face-mask-related deepfakes. To enhance the robustness of audio deepfake detection, audio-based models employ adversarial training and adaptive training techniques [22], focusing solely on the audio modality. Previous research has demonstrated the effectiveness of defending against various attack methods by utilizing projected gradient descent (PGD) as a defense mechanism, validating the proposal by Madry et al., [30]. This further supports the notion that robust defenses, specifically tailored to counter PGD attacks, can provide efficient protection against different first-order attack methods.

### 3 Adversarial Feature Similarity Learning

**Notation:** We consider  $f_\theta(\cdot)$  as the feature encoder from the deepfake detector and  $\theta$  representing its learnable model parameters. The variable  $x \in \mathcal{X}$  denotes an input frame or a video clip from a video depending on if  $f_\theta(\cdot)$  is a frame-based or a video-based deepfake detector method, where  $x$  can be either real or fake samples. In addition,  $x_{adv}$  represents the adversarial example generated using  $x \in \{x^{real}, x^{fake}\}$  where  $x^{real} \in \mathcal{X}$  and  $x^{fake} \in \mathcal{X}$  represent real and fake samples, which can be either an individual frame or a video clip depending on whether  $f_\theta(\cdot)$  is a frame-based or video-based detector.  $y \in \{0, 1\}$  denotes the label, where 0 indicates a fake class and 1 real class. Adversarial Feature Similarity Learning

#### 3.1 Overview

Our objective is to develop an adversarially robust deepfake detector that effectively mitigates the impact of adversarial attacks while preserving the performance on unperturbed data. We address this problem by discerning the features of real and fake samples and their corresponding adversarial counterparts. We



**Fig. 1.** Framework for adversarial feature similarity learning. First, we select a pair of real and deepfake samples and create adversarial perturbation for the corresponding inputs. Then, we generate the features of real, fake, and their adversarial samples. Finally, through the proposed loss function, the model can learn a better representation to separate real samples from fake ones along with their adversarial counterparts, where the backbone  $f_\theta$  is from the deepfake detector.

hypothesize that an adversarial attack will shift the features in the opposite direction, irrespective of whether the input is real or fake. However, deepfake detectors trained solely with adversarial training will not effectively learn the desirable features to distinguish between real and fake samples in the presence of adversarial attacks. To overcome this limitation, we proposed a novel loss function to effectively separate the two classes i.e. (Real and Fake) under most conditions. Figure 1 provides a comprehensive illustration of our framework. We adopt a three-step approach. Separating unperturbed (Real and Fake) in Section 3.2, using adversarial similarity loss to make the detector robust in Section 3.3, and finally similarity regularization loss to preserve the unperturbed performance in Section 3.4. Additionally, we formulate the final loss function in Section 3.5.

### 3.2 Deepfake Classification Loss

The deepfake classification loss is realized using a supervised loss that utilizes a logit-adjusted variant of binary cross-entropy (LBCE) [31], denoted as  $\mathcal{L}_{LBCE}(f_\theta(x), y)$  to tackle the potential issues of class imbalance. The supervised deepfake classification loss function aims to maximize the dissimilarity between real and fake samples, while simultaneously addressing the class imbalance in the dataset. This approach helps to refine the model’s discrimination abilities, enabling it to better differentiate between real and fake videos. The

deepfake classification loss,  $\mathcal{L}_{\text{dcl}}$ , is denoted as follows:

$$\mathcal{L}_{\text{dcl}} = \mathcal{L}_{L BCE}(f_{\theta}(x), y), \quad (1)$$

### 3.3 Adversarial Similarity Loss

Let  $f_{\theta}(x)$  and  $f_{\theta}(x_{\text{adv}})$  respectively represent the mapping from an input sample  $x \in \{x^{\text{real}}, x^{\text{fake}}\}$  and from an adversarial input sample  $x_{\text{adv}} \in \{x_{\text{adv}}^{\text{real}}, x_{\text{adv}}^{\text{fake}}\}$  to their corresponding embedding spaces. To perform adversarial training, we generate adversarial examples  $x_{\text{adv}}$  from  $x \in \{x^{\text{real}}, x^{\text{fake}}\}$  by employing the PGD adversarial attack method. The objective is to maximize the cosine similarity  $\text{sim}(f_{\theta}(x), f_{\theta}(x_{\text{adv}}))$  between  $x$  and  $x_{\text{adv}}$  to bring them closer as they represent the same class and to avoid adversarial examples from being misclassified to the other class (i.e., real to fake and fake to real). To achieve this objective while minimizing the final loss, we introduce adversarial similarity loss as follows

$$\mathcal{L}_{\text{asl}} = (1 - \text{sim}(f_{\theta}(x), f_{\theta}(x_{\text{adv}}))), \quad (2)$$

where  $\text{sim}(\cdot, \cdot)$  indicates the cosine similarity metric, and  $\mathcal{L}_{\text{asl}}$  is the adversarial similarity loss. We aim to minimize the dissimilarity using adversarial similarity loss, thereby maximizing the similarity.

### 3.4 Similarity Regularization Loss

To further enhance the performance, we enforce additional regularization to minimize the similarity between the paired real and fake samples from the corresponding real and deepfake samples. The similarity regularization loss minimizes similarity, which is computed through cosine similarity using the unperturbed samples from real and deepfake inputs. This process effectively creates separation between the two classes, ensuring that the detector’s unperturbed performance remains intact. By employing this approach, we not only improve the detector’s robustness but also preserve its unperturbed performance. Similarity regularized loss is calculated as follows:

$$\mathcal{L}_{\text{srl}} = \text{sim}(f_{\theta}(x^{\text{real}}), f_{\theta}(x^{\text{fake}})), \quad (3)$$

where  $\mathcal{L}_{\text{srl}}$  is the similarity regularization loss,  $\mathcal{L}$  is cosine similarity, and respectively  $x^{\text{real}}, x^{\text{fake}} \in \mathcal{X}$  are real and fake pair input.

### 3.5 Final Loss Function

In this section, we formulate the final loss function for minimization based on the three previously discussed components. The objective function is presented in Equation 4.

$$\mathcal{L}_{\text{afsl}} = \mathcal{L}_{\text{dcl}} + \beta_1 \mathcal{L}_{\text{asl}} + \beta_2 \mathcal{L}_{\text{srl}}, \quad (4)$$

where  $\mathcal{L}_{\text{dcl}}$  denotes the deepfake classification loss, while  $\mathcal{L}_{\text{asl}}$  and  $\mathcal{L}_{\text{srl}}$  correspond to the adversarial similarity loss and similarity regularization loss respectively.

**Table 1.** Video level AUC (%) for deepfake detectors when testing on each deepfake type of FF++ after training on the remaining three types. “No Attack” denotes that an adversarial attack is not applied while “PGD10” denotes that the Projected Gradient descent (PGD) attack is applied to the input. The types of deepfakes are DeepFakes (DF), FaceSwap (FS), Face2Face (F2F), and NeuralTextures (NT). **All the numbers of the baseline methods shown in this Table are reproduced based on the default settings of the officially released implementations, and the performance discrepancies from their papers may be due to the released versions or hyperparameters being different from the ones used in the experiments of the papers.**

Methods	No Attack				PGD10			
	DF	FS	F2F	NT	DF	FS	F2F	NT
RealForensics [15]	91.5	89.6	92.3	92.6	0.7	0.8	1.2	1.6
LipForensics [16]	90.8	87.1	92.0	91.8	2.8	2.4	1.0	2.6
FTCN [47]	91.6	90.0	91.6	92.8	1.0	1.2	1.7	2.4
Patch-based [4]	85.7	57.4	84.6	80.3	3.8	2.6	1.8	4.2
Xception [37]	83.1	50.6	81.5	76.9	1.0	4.7	0.8	0.2

To control the influence of the regularization terms, we set the scaling factors  $\beta_1$  and  $\beta_2$  to 1 and 0.1, respectively. By incorporating these components into the loss function  $\mathcal{L}_{\text{afsl}}$ , we aim to enhance the robustness of the detectors while preserving their unperturbed performance.

Unlike TRADES, our approach utilizes similarity regularization, capturing finer differences between real and fake videos and thereby enabling the extraction of more accurate features.

## 4 Experimental Description

### 4.1 Implementation Details

The faces are extracted through the utilization of face detection and alignment techniques and video clips comprising 25 frames. For the training process, the video clips are randomly cropped to dimensions  $140 \times 140$  and subsequently resized to  $112 \times 112$ . Horizontal flipping and grayscale transformation with a probability of 0.5 are applied along with random masking. For sequence-based deepfake detection, we employ the Channel-Separated Convolutional Network (CSN) [41]. We refer interested readers to [15,41] for more details about CSN. The optimization process utilizes the Adam optimizer with a learning rate of  $3 \times 10^{-4}$ . The model is trained for 150 epochs. We used RealForensics self-supervised pretrained on Lip Reading in the Wild (LRW) dataset <sup>5</sup>. This pretrained model serves as the starting point and provides the initial push to effectively capture the relevant features, thereby enhancing the model’s generalization capabilities.

<sup>5</sup> Pretrained model: <https://github.com/ahaliassos/RealForensics>

**Table 2.** AUC (%) scores for video-level detection on the FF++ dataset, containing four deepfake methods. Models train on three methods and test on the remaining method. We employ PGD5 for adversarial training and PGD10 for testing purposes.  $L_\infty$  is allowed distortion for adversarial attacks.  $L_\infty = 0$  Adversarial Feature Similarity Learning means no adversarial attack is applied.

Method	$L_\infty$	DF	FS	F2F	NT	Avg
RealForensics[15]	0	91.5	89.6	92.3	92.6	91.5
RealForensics	8/255	0.7	0.8	1.2	1.6	1.0
RealForensics + AT [30]	8/255	76.3	74.1	73.7	70.1	73.5
RealForensics + TRADES [46]	8/255	78.2	75.4	80.7	72.5	76.7
RealForensics + AFSL (Ours)	0	89.4	87.6	90.4	91.7	89.7
AFSL (Ours)	8/255	79.0	77.2	82.6	75.6	78.6
RealForensics + AFSL (Ours)	8/255	81.5	79.7	84.1	78.1	80.8

While for frame-based detection, we utilize XceptionNet [6] and MesoNet [1]<sup>6</sup>. We optimize the frame-based model using a Stochastic gradient descent (SGD) optimizer with a learning rate of  $2 \times 10^{-3}$  and the model is trained for 150 epochs with a batch size of 16. For further details about the frame-based model, interested readers are directed to [1,33]. Normalization ( $L_2$ -norm) is applied to all features in both sequence-based and frame-based detectors.

**Datasets:** FaceForensics++ (FF++), is comprised of 1,000 authentic videos and 4,000 deepfake videos. Unless specified otherwise, the mildly compressed version of the dataset (c23) was utilized. Other datasets used in the experiments are FaceShifter [25] and DeeperForensics [19], featuring different face-swapping techniques applied to FF++ real videos.

**Evaluation metrics:** We utilize accuracy and area under the receiver operating characteristic curve (AUC) metrics. For video-level assessment, we uniformly sample non-overlapping clips from a single video and average their predictions.

## 4.2 Victim Models: Deepfake Detectors

In our work, we assess the vulnerability of top-performing deepfake detectors to adversarial attacks. We employ video-based detectors, namely RealForensics [15], LipForensics [16], and FTCN [47]. In addition, we incorporate frame-by-frame based detectors, such as Patch-based [4] and Xception [37]. All the models are tested on each of the four methods using FF++ after training on the remaining three. Table 1 presents the AUC score of all detectors under two conditions: “No Attack,” where no adversarial attack is applied to the input, and “PGD10,” where an adversarial attack is applied to the input video.

As observed from the results, all the deepfake detectors demonstrate vulnerability to adversarial attacks. Our proposed loss function offers the advantage of seamless integration with most deepfake detectors. For the evaluation of

<sup>6</sup> Pretrained model: <https://github.com/paarthneekhara/AdversarialDeepFakes>

**Table 3.** Average video-level AUC (%) for adversarial attacks is computed by training the model on three methods and testing it on a fourth method. The reported values represent the average scores across the entire test dataset, encompassing all four methods, under both white-box and black-box adversarial attacks.

Methods	No Attack	PGD10	RWA [18]	CW2	SA[17]	UI [34]	RBB [18]
RealForensics [15]	<b>91.5</b>	1.1	1.4	0.0	0.0	0.0	36.8
RealForensics + AT [30]	78.4	73.5	79.5	78.3	63.6	66.9	80.5
RealForensics + TRADES [46]	84.0	76.7	76.1	78.1	67.2	69.3	82.4
AFSL (Ours)	87.3	78.3	79.1	79.7	68.5	68.7	86.1
RealForensics +AFSL (Ours)	89.8	<b>80.7</b>	<b>81.3</b>	<b>82.8</b>	<b>73.9</b>	<b>74.7</b>	<b>87.5</b>

**Table 4.** Video level AUC(%) for unseen datasets: DeeperForensics and FaceShifter under different adversarial attacks.

Methods	No Attack	PGD10	RWA [18]	CW2	SA [17]	UI [34]	RBB [18]
<b>DeeperForensics</b>							
RealForensics [15]	<b>93.6</b>	1.0	0.0	0.0	0.0	0.0	42.2
RealForensics + AT [30]	84.5	78.2	76.4	75.0	65.4	69.7	81.3
RealForensics + TRADES [46]	88.1	78.0	79.3	77.9	69.7	<b>85.6</b>	85.8
AFSL (Ours)	90.2	80.3	79.9	80.4	69.4	83.5	87.3
RealForensics +AFSL (Ours)	92.9	<b>83.6</b>	<b>81.2</b>	<b>83.5</b>	<b>72.8</b>	85.6	<b>89.1</b>
<b>FaceShifter</b>							
RealForensics [15]	<b>91.7</b>	1.0	1.0	1.0	0.0	1.0	37.7
RealForensics + AT [30]	83.6	76.2	74.0	75.8	60.7	73.1	79.9
RealForensics + TRADES [46]	87.1	79.3	77.9	78.2	<b>67.3</b>	72.8	84.6
AFSL (Ours)	88.2	79.6	78.1	81.4	65.1	74.6	84.3
RealForensics +AFSL (Ours)	89.4	<b>81.7</b>	<b>80.7</b>	<b>84.6</b>	67.1	<b>78.1</b>	<b>86.5</b>

our method, we select RealForensics from sequence-based detectors, along with XceptionNet and MesoNet from frame-by-frame detectors. As we cannot evaluate every detector, we choose the top-performing detector from each category.

We employ pretrained weights from RealForensics<sup>5</sup> and fine-tune LipForensics, FTCN<sup>7</sup>, Patch-based<sup>8</sup>, and Xception. We follow the exact instructions for pre-processing provided in their official code to replicate the results. While we do observe a decline in performance compared to the reported results, this could potentially be attributed to the absence of supplementary data. We solely report the reproduced results as we aim to improve model robustness against adversarial attacks.

### 4.3 Robust Cross-Manipulation Generalization

Most deepfake detectors typically conduct generalization experiments to assess their performance. These experiments involve training the detectors on

<sup>7</sup> <https://github.com/yinglinzheng/FTCN>

<sup>8</sup> <https://github.com/chail/patch-forensics>

**Table 5.** Frame-level Accuracy (%) of deepfake detector on FF++ dataset under different adversarial attacks using XceptionNet and MesoNet models.

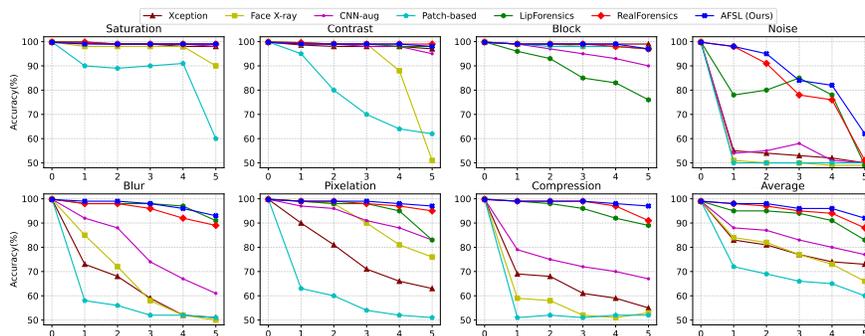
<b>XceptionNet</b>					
Methods	PGD10	CW2	RWA [18]	SA [17]	RBB [18]
Adversarial Deepfakes [18]	7.3	8.7	1.9	0.0	48.3
Vanilla AT [30]	76.0	72.7	73.1	58.9	80.1
TRADES [46]	80.8	79.4	77.2	76.8	84.7
AFSL (Ours)	<b>81.7</b>	<b>82.9</b>	<b>80.6</b>	<b>77.5</b>	<b>85.7</b>
<b>MesoNet</b>					
Adversarial Deepfakes [18]	6.2	7.3	0.0	0.0	44.6
Vanilla AT [30]	74.3	70.7	71.9	55.6	67.4
TRADES [46]	78.2	<b>77.1</b>	72.9	<b>73.7</b>	79.3
AFSL (Ours)	<b>79.5</b>	76.9	<b>74.1</b>	72.6	<b>81.3</b>

three methods and testing them on the remaining techniques using the FF++ dataset [15,38]. In this study, we follow the same protocol and introduce adversarial attacks during both the training and testing stages, utilizing a sequence-based model. To make the comparison fair, we utilized self-supervised pretrained weights of RealForensics for all methods using the CSN model. The results in Table 2 demonstrate that our proposed method achieves adversarially robust generalization to unseen adversarial deepfakes. While RealForensics achieves the best result with unperturbed input, it fails to withstand adversarial attacks. On the other hand, our proposed method performs well under both adversarially perturbed and unperturbed data compared with adversarial training (AT) and TRADES. AT is the baseline while TRADES is the state-of-the-art method in terms of both clean and robust performance. Table 3 presents the average AUC score on FF++ using white-box and black-box attacks. RBB is a black-box attack generated using ResNet3D [42], while PGD, CW, StatAttack, Universal, and RWA are white-box attacks. Our proposed method AFSL outperforms all previous defense techniques across all types of adversarial attacks.

We also evaluate the robust generalization across datasets by training the model on FF++ using all manipulation methods and testing it on two datasets, DeeperForensics and Faceshifter. Table 4 presents the AUC results for both datasets. We compare the proposed method with state-of-the-art under stronger white-box attacks and black-box attacks. The robust AUC confirms that the proposed method performs well compared to the baseline method and other defenses when exposed to various types of adversarial attacks.

#### 4.4 Evaluation on Frame-based Detectors

To showcase the effectiveness of our proposed approach, we incorporated two frame-by-frame based deepfake detectors, namely XceptionNet [6] and MesoNet [1]. These detectors are CNN-based classification models that independently classify



**Fig. 2.** Robustness to unseen distortions: Video level AUC scores (%) varying with the severity level of different distortions. Average is the mean value at each severity level.

each frame as either real or fake. Table 5 presents the performance of both models against state-of-the-art white-box adversarial attacks. For robust black-box (RBB) attacks, we generated perturbations from pre-trained clean models without accessing their parameters. This allows us to evaluate the robustness of the detectors in scenarios where they are not aware of each other internal architecture or parameters.

#### 4.5 Robustness to Common Distortions

In addition to robust generalization across different manipulations and resistance against adversarial attacks, deepfake detectors must also withstand common distortions that videos may encounter online. To assess the robustness of the detector against unfamiliar distortions, we follow the settings of [15,16]. During training on the FF++ dataset, we limit the augmentation techniques to horizontal flipping and random cropping for grayscale inputs. This approach helps prevent any potential interactions between the training distortions and those used during testing. Following [19], we use seven different types of distortions with five levels of severity. Figure 2 presents the results of each distortion with five levels of severity using the proposed method and state-of-the-art methods. The proposed method outperforms both frame-based and sequence-based methods. The inclusion of an adversarial training term in the loss function acts as regularization to increase model robustness against common distortions, which is in line with the findings of Kireev et al. [24]. In comparison to previous methods, the proposed approach performs well under most conditions, highlighting its effectiveness in tackling distortions commonly encountered in real-world scenarios.

**Table 6.** Impact of various components. Robust AUC (%) on FaceShifter (FSh) and DeeperForensics (DFo) using PGD10 adversarial attack.

Settings	Losses			AUC (%)	
	$\mathcal{L}_{dcl}$	$\mathcal{L}_{asl}$	$\mathcal{L}_{srl}$	DFo	FSh
S1	✓	✗	✗	1.0	1.0
S2	✓	✓	✗	81.3	79.1
S3	✓	✓	✓	<b>83.6</b>	<b>81.7</b>

## 5 Ablation Study

In this section, we analyze various components of our proposed method to comprehend the factors contributing to its performance. We conduct ablation experiments under PGD adversarial attack to inspect its robust generalization. The training is performed on FaceForensics++, and we evaluate the model on FaceShifter and DeeperForensics datasets, reporting the AUC score. Table 6 displays the results of the ablation study. We use deepfake classification loss as the first term in the loss function as S1 to train the model without any adversarial training. Unfortunately, this model proves to be inadequate in defending against adversarial attacks. Next, by training the model with the S2 setting without the similarity regularization, we can significantly improve the AUC scores against adversarial attacks. Finally, with all three loss components in the S3 setting, we can further improve the robustness by about 2% in AUC score.

## 6 Conclusion and Future work

This paper introduces a novel approach Adversarial Feature Similarity Learning (AFSL) for enhancing the robustness of a deepfake detector against adversarial attacks. We propose an adversarially robust loss function, specifically designed to detect fake videos even when subjected to deliberate adversarial perturbations. Our experimental results demonstrate the effectiveness of the proposed method under unperturbed input but also against common distortions. The future work will consider self-supervised learning using the proposed loss function, such as pairing real and fake samples in self-supervised adversarial defense.

## 7 ACKNOWLEDGMENT

This research is supported by National Science and Technology Council, Taiwan (R.O.C), under the grant number of NSTC-111-2634-F-002-022, 110-2221-E-001-009-MY2, 112-2634-F-001-001-MBK, and Academia Sinica under the grant number of AS-CDA-112-M09. In addition, we would like to express our gratitude for the valuable contributions and guidance from these organizations, which have been instrumental in achieving the goals of this research.

## References

1. Afchar, D., Nozick, V., Yamagishi, J., Echizen, I.: Mesonet: a compact facial video forgery detection network. In: WIFS. pp. 1–7 (2018)
2. Alnaim, N.M., Almutairi, Z.M., Alsuwat, M.S., Alalawi, H.H., Alshobaili, A., Alenezi, F.S.: Dffind: A deepfake face mask dataset for infectious disease era with deepfake detection algorithms. *IEEE Access* pp. 16711–16722 (2023)
3. Carlini, N., Wagner, D.: Adversarial examples are not easily detected: Bypassing ten detection methods. In: AIS. pp. 3–14 (2017)
4. Chai, L., Bau, D., Lim, S.N., Isola, P.: What makes fake images detectable? understanding properties that generalize. In: ECCV. pp. 103–120 (2020)
5. Chen, G., Zhao, Z., Song, F., Chen, S., Fan, L., Wang, F., Wang, J.: Towards understanding and mitigating audio adversarial examples for speaker recognition. *TDSC* (2022)
6. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: CVPR. pp. 1251–1258 (2017)
7. Deepfakes: faceswap. In: GitHub. ( Accessed: 14.06.2023) (2017), <https://github.com/deepfakes/faceswap>
8. Dong, S., Wang, J., Ji, R., Liang, J., Fan, H., Ge, Z.: Implicit identity leakage: The stumbling block to improving deepfake detection generalization. In: CVPR. pp. 3994–4004 (2023)
9. Frank, J., Eisenhofer, T., Schönherr, L., Fischer, A., Kolossa, D., Holz, T.: Leveraging frequency analysis for deep fake image recognition. In: ICML. pp. 3247–3258 (2020)
10. Gandhi, A., Jain, S.: Adversarial perturbations fool deepfake detectors. In: IJCNN. pp. 1–8 (2020)
11. Gao, G., Huang, H., Fu, C., Li, Z., He, R.: Information bottleneck disentanglement for identity swapping. In: CVPR. pp. 3404–3413 (2021)
12. Gao, Y., Wei, F., Bao, J., Gu, S., Chen, D., Wen, F., Lian, Z.: High-fidelity and arbitrary face editing. In: CVPR. pp. 16115–16124 (2021)
13. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. *ICLR* (2015)
14. Guan, J., Zhou, H., Hong, Z., Ding, E., Wang, J., Quan, C., Zhao, Y.: Delving into sequential patches for deepfake detection. *arXiv preprint arXiv:2207.02803* (2022)
15. Haliassos, A., Mira, R., Petridis, S., Pantic, M.: Leveraging real talking faces via self-supervision for robust forgery detection. In: CVPR. pp. 14950–14962 (2022)
16. Haliassos, A., Vougioukas, K., Petridis, S., Pantic, M.: Lips don’t lie: A generalisable and robust approach to face forgery detection. In: CVPR. pp. 5039–5049 (2021)
17. Hou, Y., Guo, Q., Huang, Y., Xie, X., Ma, L., Zhao, J.: Evading deepfake detectors via adversarial statistical consistency. In: CVPR. pp. 12271–12280 (2023)
18. Hussain, S., Neekhara, P., Jere, M., Koushanfar, F., McAuley, J.: Adversarial deepfakes: Evaluating vulnerability of deepfake detectors to adversarial examples. In: WACV. pp. 3348–3357 (2021)
19. Jiang, L., Li, R., Wu, W., Qian, C., Loy, C.C.: Deeperforensics-1.0: A large-scale dataset for real-world face forgery detection. In: CVPR. pp. 2889–2898 (2020)
20. Jiang, Z., Chen, T., Chen, T., Wang, Z.: Robust pre-training by adversarial contrastive learning. *NIPS* pp. 16199–16210 (2020)
21. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR. pp. 4401–4410 (2019)

22. Kawa, P., Plata, M., Syga, P.: Defense against adversarial attacks on audio deepfake detection. In: *Interspeech* (2023)
23. Khan, S., Thainimit, S., Kumazawa, I., Marukatat, S.: Text detection and recognition on traffic panel in roadside imagery. In: *2017 8th International Conference of Information and Communication Technology for Embedded Systems (IC-ICTES)*. pp. 1–6. IEEE (2017)
24. Kireev, K., Andriushchenko, M., Flammarion, N.: On the effectiveness of adversarial training against common corruptions. In: *UAI*. pp. 1012–1021 (2022)
25. Li, L., Bao, J., Yang, H., Chen, D., Wen, F.: Advancing high fidelity identity swapping for forgery detection. In: *CVPR*. pp. 5074–5083 (2020)
26. Li, Z., Yin, B., Yao, T., Guo, J., Ding, S., Chen, S., Liu, C.: Sibling-attack: Rethinking transferable adversarial attacks against face recognition. In: *CVPR*. pp. 24626–24637 (2023)
27. Liang, K., Xiao, B.: Styleless: Boosting the transferability of adversarial examples. In: *CVPR*. pp. 8163–8172 (2023)
28. Liu, B., Liu, B., Ding, M., Zhu, T., Yu, X.: Ti2net: Temporal identity inconsistency network for deepfake detection. In: *WACV*. pp. 4691–4700 (2023)
29. Lo, S.Y., Patel, V.M.: Defending against multiple and unforeseen adversarial videos. *TIP* pp. 962–973 (2021)
30. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. *ICLR* (2018)
31. Menon, A.K., Jayasumana, S., Rawat, A.S., Jain, H., Veit, A., Kumar, S.: Long-tail learning via logit adjustment. *ICLR* (2021)
32. Mumcu, F., Doshi, K., Yilmaz, Y.: Adversarial machine learning attacks against video anomaly detection systems. In: *CVPR*. pp. 206–213 (2022)
33. Neekhara, P.: Adversarialdeepfake. In: *GitHub*. ( Accessed: 14.06.2023) (2019), <https://github.com/paarthneekhara/AdversarialDeepFakes>
34. Neekhara, P., Dolhansky, B., Bitton, J., Ferrer, C.C.: Adversarial threats to deepfake detection: A practical perspective. In: *CVPR*. pp. 923–932 (2021)
35. Neekhara, P., Hussain, S., Pandey, P., Dubnov, S., McAuley, J., Koushanfar, F.: Universal adversarial perturbations for speech recognition systems. *arXiv preprint arXiv:1905.03828* (2019)
36. Qin, Y., Carlini, N., Cottrell, G., Goodfellow, I., Raffel, C.: Imperceptible, robust, and targeted adversarial examples for automatic speech recognition. In: *ICML*. pp. 5231–5240 (2019)
37. Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Nießner, M.: Faceforensics++: Learning to detect manipulated facial images. In: *CVPR*. pp. 1–11 (2019)
38. Shahzad, S.A., Hashmi, A., Khan, S., Peng, Y.T., Tsao, Y., Wang, H.M.: Lip sync matters: A novel multimodal forgery detector. In: *APSIPA*. pp. 1885–1892 (2022)
39. Songja, R., Promboot, I., Haetanurak, B., Kerdvibulvech, C.: Deepfake ai images: should deepfakes be banned in thailand? *AI and Ethics* pp. 1–13 (2023)
40. Spivak, R.: " deepfakes": The newest way to commit one of the oldest crimes. *HeinOnline* p. 339 (2018)
41. Tran, D., Wang, H., Torresani, L., Feiszli, M.: Video classification with channel-separated convolutional networks. In: *ICCV*. pp. 5552–5561 (2019)
42. Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., Paluri, M.: A closer look at spatiotemporal convolutions for action recognition. In: *CVPR*. pp. 6450–6459 (2018)

43. Wang, H., He, F., Peng, Z., Shao, T., Yang, Y.L., Zhou, K., Hogg, D.: Understanding the robustness of skeleton-based action recognition under adversarial attack. In: CVPR. pp. 14656–14665 (2021)
44. Yadlin-Segal, A., Oppenheim, Y.: Whose dystopia is it anyway? deepfakes and social media regulation. *Convergence* pp. 36–51 (2021)
45. Yang, C., Ding, L., Chen, Y., Li, H.: Defending against gan-based deepfake attacks via transformation-aware adversarial faces. In: IJCNN. pp. 1–8 (2021)
46. Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L., Jordan, M.: Theoretically principled trade-off between robustness and accuracy. ICML pp. 7472–7482 (2019)
47. Zheng, Y., Bao, J., Chen, D., Zeng, M., Wen, F.: Exploring temporal coherence for more general video face forgery detection. In: ICCV. pp. 15044–15054 (2021)