

# Towards Quantitative Evaluation of Explainable AI Methods for Deepfake Detection

Konstantinos Tsigos  
Information Technologies Institute  
CERTH, Thessaloniki, Greece  
ktsigos@iti.gr

Evlampios Apostolidis  
Information Technologies Institute  
CERTH, Thessaloniki, Greece  
apostolid@iti.gr

Spyridon Baxevanakis  
Information Technologies Institute  
CERTH, Thessaloniki, Greece  
spirosbax@iti.gr

Symeon Papadopoulos  
Information Technologies Institute  
CERTH, Thessaloniki, Greece  
papadop@iti.gr

Vasileios Mezaris  
Information Technologies Institute  
CERTH, Thessaloniki, Greece  
bmezaris@iti.gr

## ABSTRACT

In this paper we propose a new framework for evaluating the performance of explanation methods on the decisions of a deepfake detector. This framework assesses the ability of an explanation method to spot the regions of a fake image with the biggest influence on the decision of the deepfake detector, by examining the extent to which these regions can be modified through a set of adversarial attacks, in order to flip the detector's prediction or reduce its initial prediction; we anticipate a larger drop in deepfake detection accuracy and prediction, for methods that spot these regions more accurately. Based on this framework, we conduct a comparative study using a state-of-the-art model for deepfake detection that has been trained on the FaceForensics++ dataset, and five explanation methods from the literature. The findings of our quantitative and qualitative evaluations document the advanced performance of the LIME explanation method against the other compared ones, and indicate this method as the most appropriate for explaining the decisions of the utilized deepfake detector.

## CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence; Image manipulation; Visual inspection.**

## KEYWORDS

Deepfake detection, Explainable AI, Visual explanations, Evaluation framework, Adversarial image generation

## ACM Reference Format:

Konstantinos Tsigos, Evlampios Apostolidis, Spyridon Baxevanakis, Symeon Papadopoulos, and Vasileios Mezaris. 2024. Towards Quantitative Evaluation of Explainable AI Methods for Deepfake Detection. In *3rd ACM International*

**This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive version was published in Proc. MAD '24, <https://doi.org/10.1145/3643491.3660292>.**

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MAD '24, June 10, 2024, Phuket, Thailand

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-0552-6/24/06  
<https://doi.org/10.1145/3643491.3660292>

*Workshop on Multimedia AI against Disinformation (MAD '24), June 10, 2024, Phuket, Thailand.* ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3643491.3660292>

## 1 INTRODUCTION

The recent advances in the field of Generative AI have led to new and more sophisticated ways of image and video manipulation, and the creation of a new type of visual disinformation that is often referred to as deepfakes. Deepfakes are AI manipulated media in which, a person's face or body is digitally altered in an existing image or video to make them appear as someone else or to reenact them. The ongoing improvement of Generative AI technologies enables the creation of deepfakes that are increasingly difficult to detect. The latter observation, combined with the use of deepfakes for spreading disinformation, necessitates the development of effective solutions for deepfake detection. Moreover, enhancing deepfake detection methods with explanatory mechanisms would significantly improve the users' trust in these technologies and allow obtaining insights about the applied image/video manipulation procedures for creating the detected deepfake.

Despite the growing interest in building increasingly more powerful models for deepfake detection, the provision of trustworthy explanations for the output of these models has not been studied extensively. Most works on explainable deepfake detection, investigate the use of various methods that create visual explanations (usually in the form of 2D heatmaps), but evaluate the performance of methods based only on the basis of qualitative analysis over a limited set of examples [3, 18, 23, 31, 36]. Only a recent work has attempted to assess the performance of various explanation methods on two CNN-based deepfake detection models using a quantitative evaluation framework [11]. Nevertheless, their proposed framework uses explanations produced from correctly classified pristine (non-manipulated) images, in order to compare the performance of various explanation approaches. In contrast, we argue that the opposite use-case of explanations - i.e., when the model detects a deepfake - is both more meaningful and useful to the user. Moreover, their framework requires access to pairs of real-fake images, thus being non-applicable on datasets that contain only fake examples, e.g., the WildDeepfake dataset [38].

In this paper, we propose a new evaluation framework that is simpler and more widely-applicable than the one in [11]. The proposed framework takes into account the produced visual explanation for

the deepfake detector’s decision after correctly classifying a fake image, without requiring any access to its original counterpart. Based on this new framework, we evaluate the performance of five explanation methods from the literature on a state-of-the-art model for deepfake detection. Our contributions are the following:

- We explain the decisions of a state-of-the-art model for deepfake detection, that is trained to spot four different types of deepfakes, i.e., deepfake attribution.
- We perform a comparative study among five different explanation methods, aiming to identify which is the most appropriate one for the considered model.
- We propose a new evaluation framework for quantifying the ability of explanation methods to spot the most influential image regions for the decision of a deepfake detection model.

## 2 RELATED WORK

Over the last years, there is an increasing interest in the development and training of advanced network architectures for deepfake detection. However, the explanation of the decisions of these networks has been poorly investigated. In an early work, Malolan et al. [23], trained a variant of the XceptionNet [7] using a subset<sup>1</sup> of the FaceForensics++ dataset [28] and examined the use of the LIME [26] and LRP [5] methods for producing visual explanations about the outcomes of the trained model. However, the evaluation of these methods was based on a few samples and mainly focused on the robustness of the produced explanations against various affine transformations or Gaussian blurring of the input image. Xu et al. [36], utilized the representations of EfficientNet-B0 [32] and a supervised contrastive learning methodology to train a linear deepfake detector to discriminate the real from the manipulated images of the FaceForensics++ dataset [28]. In terms of explainability, Xu et al. investigated the use of the learned features only for explaining the observed detection performance, using heatmap visualizations and uniform manifold approximation and projection (UMAP). Silva et al. [31], proposed the use of an ensemble of CNNs (XceptionNet [7], EfficientNet-B3 [32]) and attention-based models for deepfake detection. They provided explanations about the regions of the images that influence the most the decision of the detector, using the Grad-CAM method [29] and focusing on the computed gradients for the attention map. Nevertheless, the produced explanations were evaluated only in a qualitative manner by taking into account only a few image samples. Jayakumar et al. [18], trained a deepfake detection model that utilizes the EfficientNet-B0 [32] as backbone and contains five dense classification layers. To produce visual explanations, they investigated the use of the Anchors [27] and LIME [26] methods, and conducted evaluations based on a limited set of examples. Aghasanli et al. [3], described a deepfake detection model that relies on Vision Transformers and can be used for distinguishing original and fake images generated by various diffusion models. For explaining the model’s output, Aghasanli et al. used SVM and xDNN [4] classifiers to understand the model’s behavior by analyzing the closest support vectors and prototypes for each classifier, respectively. The evaluation of the produced explanations though, was based on the qualitative analysis of few

samples. Haq et al. [12] described a neurosymbolic deepfake detection method that is based on the idea that deepfakes exhibit inter- or intra- modality inconsistencies in the emotional expressions of the person being manipulated. Their method performs inter- and intra-modality reasoning on emotions extracted from audio and visual modalities using a psychological and arousal valence model for deepfake detection, and provides textual explanations that localize the timestamp and identify the fake part. However, it was evaluated only in terms of deepfake and emotion detection, while its explainability dimension was discussed only theoretically. Finally, Gowrisankar et al. [11] described an evaluation framework for explanation methods, which is based on the intuition that the identified salient visual concepts by such a method after correctly classifying a real image as a non-manipulated one, could be used to flip the prediction of the detector for its fake counterpart. Initially, Gowrisankar et al., investigated the appropriateness of generic data removal/insertion approaches for modifying the spotted salient pixels or segments of the input image (e.g., zeroing, replacement with a uniform random value and blurring based on the neighboring pixels), and found out that these approaches may produce less meaningful results when applied on deepfake detection models, as they can distort facial regions and produce completely unexpected detection results (e.g., increase of deepfake detection accuracy). Based on this finding, they described a framework that applies a number of adversarial attacks (using Natural Evolution Strategies (NES) [34]) in regions of a fake image that correspond to the identified salient visual concepts after explaining the (correct) classification of its real counterpart, and evaluates the performance of an explanation method based on the observed drop in the accuracy of the deepfake detector. Thus, their evaluation framework takes the unusual step of using the produced explanation after correctly classifying a real (non-manipulated) image, in order to assess the capacity of an explanation method to explain the detection of a fake (manipulated) image.

Differently to the majority of the works described above, that evaluate explanations qualitatively using a small set of samples [3, 18, 23, 31, 36], in this work we assess the sufficiency of explanation methods to spot the regions of a manipulated image that influenced the most the deepfake detector’s decision, using a quantitative evaluation framework. Our work is most closely related with [11], but we follow a more straightforward and intuitive evaluation approach that takes into account the produced explanations for the deepfake detector’s output after correctly classifying a fake/manipulated image (which is better reflecting the task at hand). Moreover, we employ a state-of-the-art model for deepfake detection (rather than using out-of-date models, such as MesoNet [2] and XceptionNet [7]), since there is no evidence from the literature that the results for one deepfake detector can be generalized to other detectors as well.

## 3 COMPARATIVE STUDY SETUP

This section describes the employed deepfake detection model, explanation methods, and evaluation framework and measures.

<sup>1</sup>Available at: <https://github.com/ondyari/FaceForensics>

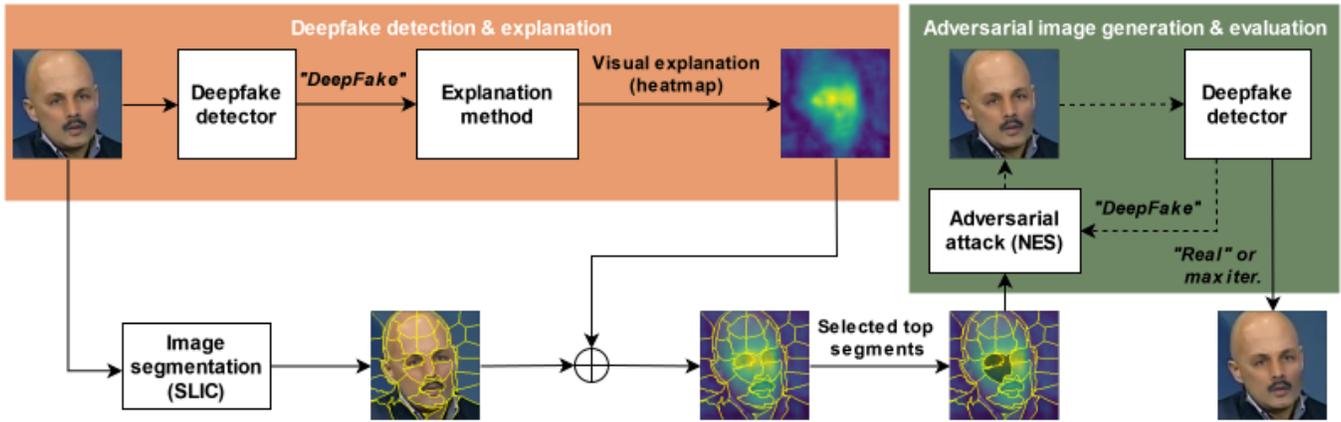


Figure 1: The processing pipeline of the proposed evaluation framework.

### 3.1 Deepfake detection model

We use a model that relies on the second version of the EfficientNet architecture [33] for deepfake detection. Building on the first version of EfficientNet - which leveraged Inverted Bottleneck convolutions (MBConv) and compound scaling to achieve high performance with fewer parameters compared to models with similar ImageNet accuracy [32] - the employed version introduces Fused Inverted Bottleneck convolutions (Fused-MBConv), leading to even faster training and improved efficiency [33]. We chose EfficientNet due to its widespread adoption, efficiency, and effectiveness in deepfake detection tasks, either as a part of an ensemble or as a backbone of more advanced methods [30, 37]. Notably, an ensemble of five EfficientNet-B7 models achieved the winning performance in Meta’s DFDC challenge [9]. Moreover, EfficientNet has been shown to outperform alternative CNN architectures, such as XceptionNet [7] and MesoNet [2] (that were taken into account in [11]), on various deepfake datasets [15, 19, 24]. Finally, it achieves similar performance to other vanilla CNNs on the ForgeryNet dataset [13] while requiring fewer parameters.

### 3.2 Explanation methods

We produce visual explanations by highlighting the regions of the image (or video frame) with the biggest influence on the deepfake detection model’s decision. As depicted in the orange coloured part of Fig. 1, we explain the outcome of the deepfake detector for a given fake image, using 2D heatmaps that represent the significance of different parts of the input image using a color scale. In our study, we consider the following explanation methods:

- **Grad-CAM++** [6], is a back-propagation-based method that generates visual explanations by leveraging the information flow (gradients) during the back-propagation process. It extends the Grad-CAM method [29], by calculating a weighted combination of the positive partial derivatives of the last convolutional layer with respect to a specific class score in order to generate the visual explanation. In this way, it provides better (more complete) object localization and is capable of explaining occurrences of multiple instances of a given object in a single image.

- **RISE** [25], is a perturbation-based method that produces visual explanations by randomly masking out portions of the input image and assessing their impact on the model’s output. Initially, this method generates a set of binary masks that are used to occlude regions of the input image and produce a set of perturbed images. Then, it feeds these perturbed images to the model, gets the model’s predictions for each one of them and uses them to weight the corresponding binary masks. Finally, it creates the visual explanation by aggregating the weighted masks together.
- **SHAP** [22], is an attribution-based method that leverages the Shapley values from game theory. It constructs an additive feature attribution model that attributes an effect to each input feature and sums the effects, i.e., SHAP values, as a local approximation of the output. More specifically, Shapley values assign importance scores to the individual pixels of the input image by treating them as players in a coalition game, with each player’s presence or absence affecting the final outcome. The payout of the grand coalition is the prediction, or in our case the explanation, and Shapley values are used to divide this payout equally among pixels, by leveraging the model predictions for the perturbed images, to assess the contribution of each pixel to the prediction.
- **LIME** [26], is a perturbation-based method that creates visual explanations by randomly masking out portions of an input image to assess their impact on the model’s output. The fundamental idea behind LIME is the approximation of the model’s behavior locally (i.e. around a specific instance) by generating a simpler, interpretable model. To this end, the input image is initially segmented and perturbed by randomly masking segments of it. Then, the perturbed images are given to the model that outputs its predictions. Finally, using a linear model (e.g., a linear regressor), LIME fits the binary masks of each perturbation to the corresponding scores and constructs the visual explanation by examining the coefficients/weights that emerge from this simpler model.
- **SOBOL** [10], is an attribution-based method that employs a mathematical concept called Sobol’ indices (after Ilya M.



**Figure 2: The produced explanations by the LIME method (the best performing one according to the results in Section 4), for three non-manipulated images of the FaceForensics++ dataset, that were correctly classified as “real”.**

Sobol’<sup>2</sup>), to identify the contribution of the input variables on the variance of the model’s output. Using a Quasi-Monte Carlo sequence, SOBOL generates a set of real-valued masks, which are then applied to an input image using perturbation functions such as blurring, to generate the perturbations. The resulting images are then forwarded to the model to get the prediction scores. By analyzing the relationship between the masks and their associated prediction scores, SOBOL estimates the total order of Sobol’ indices and creates a visual explanation by highlighting the importance of each region.

### 3.3 Evaluation framework and measures

Based on the reported findings in [11], about the adequacy of generic data removal/insertion approaches for perturbing the input image, we also do not apply such approaches on the image regions that have been promoted by an explanation method, in order to assess this method’s performance. We evaluate the performance of an explanation method by extending the evaluation framework in [11] so that it takes into account the produced explanations for fake images. We argue that the provision of an explanation after detecting a fake image is more meaningful for the user, as it can give clues about regions of the image (the highlighted ones by the visual explanation) that were found to be manipulated. On the contrary, the provided explanation after classifying an image as “real” would demarcate specific regions of the image as non-manipulated (see Fig. 2), while someone would expect that the entire image has not been manipulated at all.

Let us assume a fake image and the produced visual explanation for the deepfake detector’s decision, by an explanation method (see the orange coloured part of Fig. 1). We assess the performance of this method by examining the extent to which the indicated regions in the visual explanation as the most important ones, can be used to flip the deepfake detector’s decision (and thus classify the image as “real”). For this, we segment the input image into super-pixel segments using the SLIC algorithm [1]. Then, we quantify the contribution of each segment to the deepfake detector’s decision by overlaying the created visual explanation to the segmented image and averaging the scores of the explanation for the pixels of the segment - as a note, in the case of LIME [26] we pass the SLIC-based segmentation mask of the input image and get the

---

#### Algorithm 1 Adversarial image generation

---

**Parameters:** Search variance  $\sigma$ , Number of samples  $n$ , Image dimension  $N$ , Maximum number of iterations  $itr$ , Maximum distortion  $\epsilon$ , Learning rate  $\alpha$

**Input:** Deepfake image  $x$ , Deepfake detector  $F$ , Binary mask  $M$

**Output:** Adversarial image  $x_{adv}$

```

1:  $x_{adv} = x$ ,
2: for  $i = 1 \rightarrow itr$  do
3:   if  $F(x_{adv}) = real$  then
4:     return  $x_{adv}$ 
5:    $g = 0$ 
6:   for  $j = 1 \rightarrow n$  do
7:      $u_j = N(0_N, I_{N,N})$ 
8:      $g = g + F(x_{adv}[M] + \sigma u_j[M])_{label=real} \cdot u_j[M]$ 
9:      $g = g - F(x_{adv}[M] + \sigma u_j[M])_{label=real} \cdot u_j[M]$ 
10:     $g = \frac{1}{2n\sigma}g$ 
11:     $x_{adv}[M] = x_{adv}[M] + clip_{\epsilon}(\alpha \cdot sign(g))$ 
12: return  $x_{adv}$ 

```

---

top-k scoring segments directly. Following, we focus on the top-k scoring segments and apply NES to progressively generate a variant of the input image that is classified as “real” by the deepfake detector. This iterative process, called adversarial image generation and evaluation, is illustrated in the green coloured area of Fig. 1. In each step of this process, we produce a variant of the input image by adding noise to the regions corresponding to the top-k scoring segments, following the steps in Alg. 1. The adversarial image generation and evaluation process stops if the deepfake detector classifies the adversarial image as “real” or a maximum number of iterations is reached.

To quantify the performance of an explanation method, we calculate the accuracy of the deepfake detection model on the set of returned adversarial images after the completion of the adversarial image generation and evaluation process, when the adversarial attacks target the top-1, top-2 and top-3 scoring segments of the input images by the method. This measure ranges in  $[0, 1]$ , where the upper boundary denotes a 100% detection accuracy. We anticipate a larger decrease in accuracy for explanation methods that spot the most influential regions of the input image for the deepfake detector’s decision, more effectively. Complementary to the aforementioned measure, we quantify the sufficiency of explanation methods to spot the most influential image regions for the deepfake detector, by calculating also the difference in the detector’s output after applying adversarial attacks to the top-1, top-2 and top-3 scoring segments (following the paradigm in [20]). This measure ranges in  $[0, 1]$ , where low/high sufficiency scores indicate that the top-k scoring segments by the explanation method have low/high impact to the deepfake detector’s decision, and thus the produced visual explanation exhibits low/high sufficiency.

## 4 EXPERIMENTS

This section discusses the utilized dataset and implementation details, and reports the findings of our quantitative and qualitative evaluations.

<sup>2</sup>[https://en.wikipedia.org/wiki/Ilya\\_M.\\_Sobol%27](https://en.wikipedia.org/wiki/Ilya_M._Sobol%27)

**Table 1: The accuracy of the employed deepfake detection model for the different types of fakes in the FaceForensics++ dataset, on the original set of images (second row) and the adversarially-generated variants of them after modifying the image regions corresponding to the top-1, top-2 and top-3 scoring segments based on the different explanation methods. Best scores in bold and second best scores underlined.**

	DF			F2F			FS			NT		
Original Accuracy	0.978			0.977			0.982			0.924		
	Top 1	Top 2	Top 3	Top 1	Top 2	Top 3	Top 1	Top 2	Top 3	Top 1	Top 2	Top 3
Grad-CAM++	0.781	0.644	0.571	0.864	0.798	0.737	0.887	0.808	0.728	0.601	0.481	0.432
RISE	0.877	0.766	0.686	0.843	0.710	0.622	0.896	0.809	0.734	0.783	0.637	0.513
SHAP	0.813	0.609	<u>0.450</u>	0.846	0.739	0.637	0.876	<u>0.702</u>	<b>0.543</b>	0.686	0.497	0.344
LIME	<b>0.735</b>	<b>0.440</b>	<u>0.245</u>	<b>0.803</b>	<b>0.633</b>	<b>0.484</b>	<b>0.864</b>	<u>0.698</u>	<u>0.559</u>	<b>0.579</b>	<b>0.340</b>	<b>0.197</b>
SOBOL	<u>0.750</u>	<u>0.591</u>	0.490	<u>0.816</u>	<u>0.653</u>	<u>0.512</u>	<u>0.874</u>	0.703	0.574	<u>0.621</u>	<u>0.417</u>	<u>0.313</u>

#### 4.1 Dataset and implementation details

Our experiments were conducted on the FaceForensics++ dataset [28]. This dataset contains 1000 original videos and 4000 fake videos created using one of the following four classes of AI-based manipulation (1000 videos per class): “FaceSwap” (FS), “DeepFakes” (DF), “Face2Face” (F2F), and “NeuralTextures” (NT). The videos of the FS class were created via a graphics-based approach that transfers the face region from a source to a target video. The videos of the DF class were produced using autoencoders to replace a face in a target sequence with a face in a source video or image collection. The videos of the F2F class were obtained by a facial reenactment system that transfers the expressions of a source to a target video while maintaining the identity of the target person. The videos of the NT class were generated by modifying the facial expressions corresponding to the mouth region, using a patch-based GAN-loss as utilized in Pix2Pix [17]. The dataset is divided into training, validation, and test sets, comprised of 720, 140 and 140 videos, respectively.

For deepfake detection, we sampled the videos keeping 1 frame per second and used the RetinaFace face detector [8] to obtain bounding boxes for the present faces. Following suggestions in [28], we enlarged each bounding box by a factor of 1.3 to capture relevant background information that might aid in discriminating between real and fake samples. The cropped faces were stored and used as input to train and test the deepfake detector. For training, we leveraged a pre-trained model on the ImageNet 1K dataset obtained from the timm library [35]. Then, the deepfake detection model was trained for 30 epochs using the AdamW optimizer [21] with a learning rate of  $5 \times 10^{-5}$  and a weight decay of  $1 \times 10^{-1}$ , and the Cross-Entropy loss for multiclass classification. To mitigate overfitting and improve generalization, we employed the following data augmentation techniques: Random Erasing, Random Resized Crop, and AugMix [14]. Additionally, to improve robustness to unseen data and encourage the model to learn more reliable features, we incorporated Stochastic Depth [16] with a drop path rate of  $4 \times 10^{-1}$ . As a result, there was a 40% chance of dropping a residual block connection during each forward pass.

To obtain the data for evaluating the different explanation methods we followed the approach in [11]. In particular, we used 127 videos from each different class of the test set and we sampled 10 frames per video, thus creating four sets of 1270 images. The generation of visual explanations was based on the following settings:

- For **Grad-CAM++**, we took the average of all convolutional 2D layers.
- For **RISE**, we set the number of masks equal to 4000 and kept all the other parameters with their default values.
- For **SHAP**, we set the number of evaluations equal to 2000 and used a blurring mask with kernel size equal to 128.
- For **LIME**, we set the number of perturbations equal to 2000 and used the SLIC algorithm with a target number of segments equal to 50.
- For **SOBOL**, we set the grid size equal to 8 and the number of design equal to 32, and kept all the other parameters with their default values.

With respect to NES, we set: the number of maximum iterations equal to 50, the learning rate equal to  $1/255$ , the maximum distortion equal to  $16/255$ , the search variance equal to 0.001, and the number of samples equal to 40. All experiments were carried out on NVIDIA RTX 4090 GPU cards. The code for reproducing the reported results is publicly-available at: <https://github.com/IDT-ITI/XAI-Deepfakes>

#### 4.2 Quantitative results

Table 1 reports the accuracy of the employed deepfake detection model for the different types of fakes in the FaceForensics++ dataset, on the original set of images (second row) and the adversarially-generated variants of them after modifying the image regions corresponding to the top-1, top-2 and top-3 scoring segments according to the different explanation methods. As shown in this table, the used deepfake detection model exhibits very high performance on all types of fakes of this dataset (achieving approx. 98% accuracy on DF, F2F and FS and over 92% on NT), documenting its state-of-the-art performance. With respect to the considered explanation methods, LIME appears to be the most effective one, as it is associated with the largest decrease in the detection accuracy for all types of fakes and in almost all experimental settings. As expected, the observed accuracy decrease is smaller when the adversarial image is generated based on the top-1 scoring segment and significantly larger when the adversarial attack is performed on the top-2 and top-3 scoring segments. However, this decrease is even more pronounced in the case of LIME. Therefore, LIME appears to be more effective compared to the other methods at highlighting the most influential segment of the input image for the decisions of the used deepfake detector, and noticeably better at spotting the top-2 or top-3 image segments with the highest impact on the

**Table 2: The sufficiency scores of the considered explanation methods for the different types of fakes in the FaceForensics++ dataset, after modifying the top-1, top-2 and top-3 scoring segments of the input images. Best scores in bold and second best scores underlined.**

	DF			F2F			FS			NT		
	Top 1	Top 2	Top 3	Top 1	Top 2	Top 3	Top 1	Top 2	Top 3	Top 1	Top 2	Top 3
Grad-CAM++	0.148	0.253	0.310	0.069	0.115	0.162	0.063	0.113	0.160	0.194	0.251	0.280
RISE	0.087	0.162	0.219	0.091	0.173	0.223	0.060	0.114	0.157	0.115	0.204	0.273
SHAP	0.137	0.300	<u>0.402</u>	0.092	0.158	0.222	0.073	<u>0.181</u>	<b>0.269</b>	0.167	0.282	0.357
LIME	<b>0.195</b>	<b>0.408</b>	<b>0.539</b>	<b>0.121</b>	<b>0.238</b>	<b>0.334</b>	<b>0.087</b>	<b>0.189</b>	<u>0.262</u>	<b>0.233</b>	<b>0.363</b>	<b>0.431</b>
SOBOL	<u>0.166</u>	<u>0.277</u>	0.352	<u>0.108</u>	<u>0.212</u>	<u>0.296</u>	<u>0.078</u>	0.180	0.259	<u>0.198</u>	<u>0.302</u>	<u>0.362</u>

**Table 3: Comparison of the obtained deepfake detection accuracy scores using our evaluation framework and the framework proposed in [11]. Best scores in bold and second best scores underlined.**

	Our framework				Framework in [11]			
	DF	F2F	FS	NT	DF	F2F	FS	NT
Original Accuracy	0.930	1.000	1.000	0.973	0.930	1.000	1.000	0.973
Grad-CAM++	0.462	0.730	0.823	0.544	0.329	0.486	0.605	0.456
RISE	0.538	0.527	0.714	0.544	0.285	0.459	0.585	0.456
SHAP	<u>0.177</u>	0.547	<u>0.558</u>	0.415	<u>0.247</u>	<u>0.378</u>	<u>0.537</u>	<b>0.286</b>
LIME	<b>0.101</b>	<b>0.250</b>	<b>0.476</b>	<b>0.265</b>	0.367	0.473	0.599	<u>0.299</u>
SOBOL	0.335	<u>0.324</u>	0.571	<u>0.388</u>	<b>0.222</b>	<b>0.331</b>	<b>0.517</b>	0.354

**Table 4: Comparison of the obtained sufficiency scores using our evaluation framework and the framework proposed in [11]. Best scores in bold and second best scores underlined.**

	Our framework				Framework in [11]			
	DF	F2F	FS	NT	DF	F2F	FS	NT
Grad-CAM++	0.355	0.211	0.116	0.232	0.397	0.350	0.252	0.319
RISE	0.280	0.310	0.173	0.276	<u>0.447</u>	0.364	0.262	0.317
SHAP	<u>0.483</u>	0.296	<u>0.276</u>	<u>0.356</u>	0.427	<u>0.382</u>	<u>0.279</u>	<u>0.411</u>
LIME	<b>0.565</b>	<b>0.495</b>	<b>0.323</b>	<b>0.441</b>	0.370	0.331	0.246	<b>0.414</b>
SOBOL	0.431	<u>0.434</u>	<u>0.276</u>	0.345	<b>0.493</b>	<b>0.444</b>	<b>0.300</b>	0.376

detector’s decision. Concerning the remaining methods, SOBOL seems to be the most competitive in most cases, while SHAP shows good performance in the case of DF and FS samples when spotting the top-2 or top-3 regions of the image. Finally, a comparison of the reported results across the different types of fakes, reveals that the different explanation methods can more effectively explain the detection of DF and NT classes, while the explanation of fakes from the remaining two classes is a more challenging task.

Table 2 presents the sufficiency scores of the considered explanation methods for the different types of fakes in the FaceForensics++ dataset, after performing adversarial attacks at the top-1, top-2 and top-3 scoring segments of the input images. These scores seem to be aligned with the results in Table 1, demonstrating once again, that LIME performs consistently good for all the considered types of fakes and numbers of top-scoring segments. Moreover, its effectiveness in spotting the most influential regions of the images is more pronounced when taking into account the top-3 scoring segments according to the produced visual explanation. As before, SOBOL is the second best method and SHAP performs comparatively good in specific occasions. Finally, the most challenging cases in terms

of visual explanation, still remain the ones associated with fakes of the F2F and FS classes.

Finally, we compared the obtained results after applying the proposed evaluation framework and the one in [11]. As a note, this comparison was based on a subset of (randomly) selected images (150 per class of fakes) to limit the computational needs of the experiment. The scores about the deepfake detector’s accuracy and the explanation method’s sufficiency are reported in Tables 3 and 4, respectively. These demonstrate that the different frameworks lead to different outcomes about the performance and the ranking of the considered explanation methods. Once again, LIME is the best-performing method for the selected subset of images according to our framework, while the framework from [11] points to SOBOL as the most effective method. The observed difference is explained by the fact that the two frameworks base their evaluations on different conditions. The framework from [11] assesses the performance of an explanation method by taking into account the produced explanations for correctly classified non-manipulated images. On the contrary, our framework focuses on the produced explanations for manipulated images that were classified as deepfakes, since

highlighting the image regions that were perceived as manipulated is more meaningful for the users.

### 4.3 Qualitative results

The top row of Fig. 3 shows four different images (sampled video frames) of the FaceForensics++ dataset and the next row contains their AI manipulated variants, where each variant is associated with a different type of manipulation. The remaining rows present the produced visual explanations by the examined methods. As illustrated in these rows, LIME successfully spots: i) the regions close to the eyes and mouth that have been modified in the case of the DF sample, ii) the regions around the nose and the cheeks that have been changed in the case of the F2F sample, iii) the regions close to the left eye and cheek that have been altered in the case of the FS sample, and iv) the regions close to the mouth and chin that have been manipulated in the case of the NT sample. With respect to the other explanation methods, Grad-CAM++ correctly focuses on regions close to the eyes in the DF and FS samples and close to the chin in the case of the NT sample. However, it fails to clearly indicate regions in the case of the F2F sample and to spot manipulations around the mouth in the case of the DF sample. RISE seems to produce explanations that highlight irrelevant (see the F2F and FS samples) or non-manipulated regions (see the NT sample) of the image, while also failing to spot the manipulated ones (see the DF sample). Finally, SHAP and SOBOL appear to perform well compared with LIME, as in most cases, they provide explanations that indicate the altered regions of the images. This finding is aligned with the performance of these methods according to the conducted quantitative evaluation.

## 5 CONCLUSIONS

In this paper, we presented a new evaluation framework for explainable AI methods for deepfake detection, which measures the capacity of such methods to spot the most influential regions of the input image through an adversarial image generation and evaluation process that aims to flip the detector's decision. We applied this framework on a state-of-the-art model for deepfake detection and five explanation methods from the literature. Our experimental results demonstrate the competitive performance of the LIME explanation method across all different types of fakes, and its competency to produce meaningful explanations for the employed deepfake detection model.

## ACKNOWLEDGMENTS

This work was supported by the EU Horizon Europe and Horizon 2020 programmes under grant agreements 101070190 AI4TRUST and 951911 AI4Media, respectively.

## REFERENCES

- [1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. 2012. SLIC Superpixels Compared to State-of-the-Art Superpixel Methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34, 11 (2012), 2274–2282. <https://doi.org/10.1109/TPAMI.2012.120>
- [2] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. 2018. MesoNet: a Compact Facial Video Forgery Detection Network. In *2018 IEEE International Workshop on Information Forensics and Security, WIFS 2018, Hong Kong, China, December 11-13, 2018*. IEEE, 1–7. <https://doi.org/10.1109/WIFS.2018.8630761>
- [3] Agil Aghasani, Dmitry Kangin, and Plamen Angelov. 2023. Interpretable-through-prototypes deepfake detection for diffusion models. In *2023 IEEE/CVF Int. Conf. on Computer Vision Workshops (ICCVW)*. IEEE Computer Society, Los Alamitos, CA, USA, 467–474. <https://doi.org/10.1109/ICCVW60793.2023.00053>
- [4] Plamen Angelov and Eduardo Soares. 2020. Towards explainable deep neural networks (xDNN). *Neural Networks* 130 (2020), 185–194. <https://doi.org/10.1016/j.neunet.2020.07.010>
- [5] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation. *PLOS ONE* 10, 7 (07 2015), 1–46. <https://doi.org/10.1371/journal.pone.0130140>
- [6] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. 2018. Grad-CAM++: Generalized Gradient-Based Visual Explanations for Deep Convolutional Networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. 839–847. <https://doi.org/10.1109/WACV.2018.00097>
- [7] Francois Chollet. 2017. Xception: Deep Learning with Depthwise Separable Convolutions. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, Los Alamitos, CA, USA, 1800–1807. <https://doi.org/10.1109/CVPR.2017.195>
- [8] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. 2020. RetinaFace: Single-Shot Multi-Level Face Localisation in the Wild. In *2020 IEEE/CVF Conf. on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. Computer Vision Foundation / IEEE, 5202–5211. <https://doi.org/10.1109/CVPR42600.2020.00525>
- [9] Brian Dolhansky, Russ Howes, Ben Pfaff, Nicole Baram, and Cristian Canton-Ferrer. 2019. The Deepfake Detection Challenge (DFDC) Preview Dataset. *CoRR abs/1910.08854* (2019). arXiv:1910.08854 <http://arxiv.org/abs/1910.08854>
- [10] Thomas Fel, Remi Cadene, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas Serre. 2021. Look at the Variance! Efficient Black-box Explanations with Sobol-based Sensitivity Analysis. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [11] Balachandar Gowrisankar and Vrizzlynn L.L. Thing. 2024. An adversarial attack approach for explainable AI evaluation on deepfake detection models. *Computers & Security* 139 (2024), 103684. <https://doi.org/10.1016/j.cose.2023.103684>
- [12] Ijaz Ul Haq, Khalid Mahmood Malik, and Khan Muhammad. 2023. Multimodal Neurosymbolic Approach for Explainable Deepfake Detection. *ACM Trans. Multimedia Comput. Commun. Appl.* (sep 2023). <https://doi.org/10.1145/3624748> Just Accepted.
- [13] Yanan He, Bei Gan, Siyu Chen, Yichun Zhou, Guojun Yin, Luchuan Song, Lu Sheng, Jing Shao, and Ziwei Liu. 2021. ForgeryNet: A Versatile Benchmark for Comprehensive Forgery Analysis. In *IEEE Conf. on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 4360–4369. <https://doi.org/10.1109/CVPR46437.2021.00434>
- [14] Dan Hendrycks, Norman Mu, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. 2020. AugMix: A Simple Data Processing Method to Improve Robustness and Uncertainty. In *8th Int. Conf. on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. <https://openreview.net/forum?id=S1gmrxHFvB>
- [15] Baojin Huang, Zhongyuan Wang, Jifan Yang, Jiaxin Ai, Qin Zou, Qian Wang, and Dengpan Ye. 2023. Implicit Identity Driven Deepfake Face Swapping Detection. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*. IEEE, 4490–4499. <https://doi.org/10.1109/CVPR52729.2023.00436>
- [16] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q. Weinberger. 2016. Deep Networks with Stochastic Depth. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV (Lecture Notes in Computer Science, Vol. 9908)*, Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (Eds.). Springer, 646–661. [https://doi.org/10.1007/978-3-319-46493-0\\_39](https://doi.org/10.1007/978-3-319-46493-0_39)
- [17] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. 2017. Image-to-Image Translation with Conditional Adversarial Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, Los Alamitos, CA, USA, 5967–5976. <https://doi.org/10.1109/CVPR.2017.632>
- [18] Krishnakripa Jayakumar and Nimalaprakas Skandhakumar. 2022. A Visually Interpretable Forensic Deepfake Detection Tool Using Anchors. In *2022 7th Int. Conf. on Information Technology Research (ICITR)*. 1–6. <https://doi.org/10.1109/ICITR57877.2022.9993294>
- [19] Xin Li, Rongrong Ni, Pengpeng Yang, Zhiqiang Fu, and Yao Zhao. 2023. Artifacts-Disentangled Adversarial Learning for Deepfake Detection. *IEEE Trans. Circuits Syst. Video Technol.* 33, 4 (2023), 1658–1670. <https://doi.org/10.1109/TCSVT.2022.3217950>
- [20] Yibing Liu, Haoliang Li, Yangyang Guo, Chenqi Kong, Jing Li, and Shiqi Wang. 2022. Rethinking Attention-Model Explainability through Faithfulness Violation Test. In *Proc. of the 39th Int. Conf. on Machine Learning (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (Eds.). PMLR, 13807–13824. <https://proceedings.mlrpress/v162/liu22i.html>

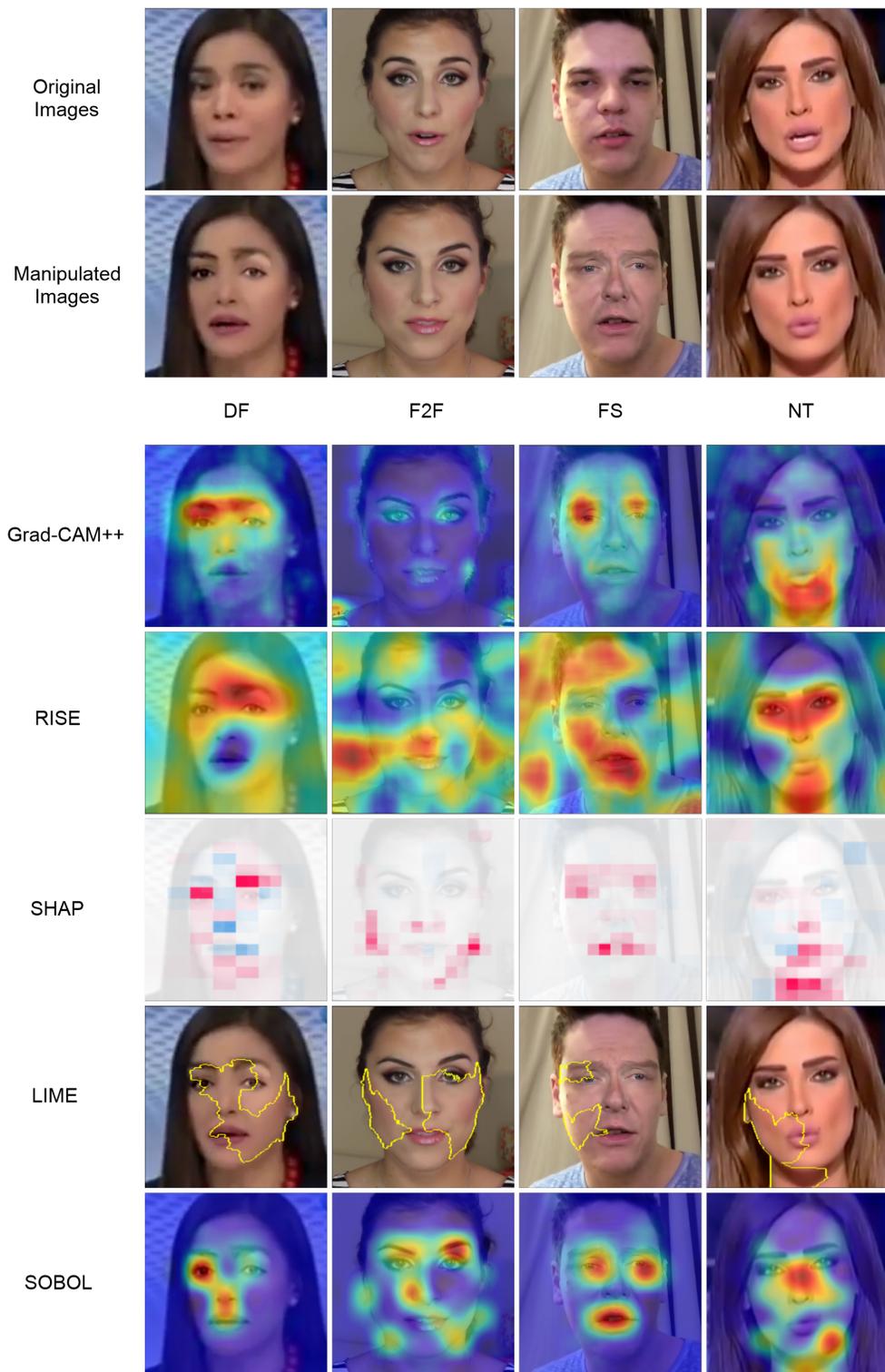


Figure 3: The obtained visual explanations from the considered explanation methods for four different images of the Face-Forensics++ dataset (one per different type of manipulation). In terms of visualization, we adopt the default supported format by each explanation method.

- [21] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *7th Int. Conf. on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. <https://openreview.net/forum?id=Bkg6RiCqY7>
- [22] Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Proc. of the 31st Int. Conf. on Neural Information Processing Systems (Long Beach, California, USA) (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 4768–4777.
- [23] Badhrinarayan Malolan, Ankit Parekh, and Faruk Kazi. 2020. Explainable Deep-Fake Detection Using Visual Interpretability Methods. In *2020 3rd Int. Conf. on Information and Computer Technologies (ICICT)*. 289–293. <https://doi.org/10.1109/ICICT50521.2020.00051>
- [24] Aakash Varma Nadimpalli and Ajita Rattani. 2023. Facial Forgery-based Deepfake Detection using Fine-Grained Features. *CoRR* abs/2310.07028 (2023). <https://doi.org/10.48550/ARXIV.2310.07028> arXiv:2310.07028
- [25] Vitali Petsiuk, Abir Das, and Kate Saenko. 2018. RISE: Randomized Input Sampling for Explanation of Black-box Models. *CoRR* abs/1806.07421 (2018). arXiv:1806.07421 <http://arxiv.org/abs/1806.07421>
- [26] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *Proc. of the 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (San Francisco, California, USA) (KDD '16)*. Association for Computing Machinery, New York, NY, USA, 1135–1144. <https://doi.org/10.1145/2939672.2939778>
- [27] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: high-precision model-agnostic explanations. In *Proc. of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence (New Orleans, Louisiana, USA) (AAAI'18/IAAI'18/EAAI'18)*. AAAI Press, Article 187, 9 pages.
- [28] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Niessner. 2019. FaceForensics++: Learning to Detect Manipulated Facial Images. In *2019 IEEE/CVF Int. Conf. on Computer Vision (ICCV)*. IEEE Computer Society, Los Alamitos, CA, USA, 1–11. <https://doi.org/10.1109/ICCV.2019.00009>
- [29] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. In *2017 IEEE Int. Conf. on Computer Vision (ICCV)*. 618–626. <https://doi.org/10.1109/ICCV.2017.74>
- [30] Kaede Shiohara and Toshihiko Yamasaki. 2022. Detecting Deepfakes with Self-Blended Images. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 18699–18708. <https://doi.org/10.1109/CVPR52688.2022.01816>
- [31] Samuel Henrique Silva, Mazal Bethany, Alexis Megan Votto, Ian Henry Scarf, Nicole Beebe, and Peyman Najafirad. 2022. Deepfake forensics analysis: An explainable hierarchical ensemble of weakly supervised models. *Forensic Science International: Synergy* 4 (2022), 100217.
- [32] Mingxing Tan and Quoc V. Le. 2019. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. In *Proc. of the 36th Int. Conf. on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 6105–6114. <http://proceedings.mlr.press/v97/tan19a.html>
- [33] Mingxing Tan and Quoc V. Le. 2021. EfficientNetV2: Smaller Models and Faster Training. In *Proc. of the 38th Int. Conf. on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 10096–10106. <http://proceedings.mlr.press/v139/tan21a.html>
- [34] Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. 2014. Natural Evolution Strategies. *Journal of Machine Learning Research* 15, 27 (2014), 949–980. <http://jmlr.org/papers/v15/wierstra14a.html>
- [35] Ross Wightman. 2019. PyTorch Image Models. <https://github.com/rwightman/pytorch-image-models>. <https://doi.org/10.5281/zenodo.4414861>
- [36] Ying Xu, Kiran Raja, and Marius Pedersen. 2022. Supervised Contrastive Learning for Generalizable and Explainable DeepFakes Detection. In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*. 379–389. <https://doi.org/10.1109/WACVW54805.2022.00044>
- [37] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. 2021. Multi-Attentional Deepfake Detection. In *IEEE Conf. on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 2185–2194. <https://doi.org/10.1109/CVPR46437.2021.00222>
- [38] Bojia Zi, Minghao Chang, Jingjing Chen, Xingjun Ma, and Yu-Gang Jiang. 2020. WildDeepfake: A Challenging Real-World Dataset for Deepfake Detection. In *Proc. of the 28th ACM Int. Conf. on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 2382–2390. <https://doi.org/10.1145/3394171.3413769>