

Identification of Deep Network Generated Images Using Disparities in Color Components

Haodong Li^{a,b}, Bin Li^{a,b,*}, Shunquan Tan^{a,b}, Jiwu Huang^{a,b}

^a *Guangdong Key Laboratory of Intelligent Information Processing, Shenzhen Key Laboratory of Media Security, Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen University, Shenzhen 518060, China*

^b *Shenzhen Institute of Artificial Intelligence and Robotics for Society, Shenzhen 518172, China*

Abstract

With the powerful deep network architectures, such as generative adversarial networks, one can easily generate photorealistic images. Although the generated images are not dedicated for fooling human or deceiving biometric authentication systems, research communities and public media have shown great concerns on the security issues caused by these images. This paper addresses the problem of identifying deep network generated (DNG) images. Taking the differences between camera imaging and DNG image generation into considerations, we analyze the disparities between DNG images and real images in different color components. We observe that the DNG images are more distinguishable from real ones in the chrominance components, especially in the residual domain. Based on these observations, we propose a feature set to capture color image statistics for identifying DNG images. Additionally, we evaluate several detection situations, including the training-testing data are matched or mismatched in image sources or generative models and detection with only real images. Extensive experimental results show that the proposed method can accurately identify DNG images and outperforms existing methods when the training and testing data are mismatched. Moreover, when the GAN model is unknown, our methods also achieves good performance with one-class classification by using only real images for training.

Keywords: Image generative model, generative adversarial networks, fake image identification, color disparities, statistical feature.

*Corresponding author

Email addresses: lihaodong@szu.edu.cn (Haodong Li), libin@szu.edu.cn (Bin Li), tansq@szu.edu.cn (Shunquan Tan), jwhuang@szu.edu.cn (Jiwu Huang)

1. Introduction

In recent years, we have witnessed the inspiring development of image generative models [1, 2, 3]. Traditionally, image generative models could only create simple image textures, and the image contents were far from realistic. Therefore, it was not difficult to differentiate between such generated images and real images with naked eyes. However, the situation has changed with the tremendous advancement of machine learning. With the modern deep network architectures, especially the generative adversarial networks (GANs) [1], the quality of generated images has been dramatically improved [4, 5, 6], and thus it is no longer easy to identify the deep network generated (DNG) images with human visual system. Although the advancements of image generative models have facilitated many applications, such as image super-resolution [7, 8], image translation [9, 10], and image inpainting [11], the photorealistic DNG images may also lead to many serious security risks. For example, the generated scenes can be used to falsify images or videos and fabricate fake news, the generated faces can be posted on social networks to counterfeit personal information or be used to attack biometric authentication systems. Recently, both public media [12] and research communities [13] have shown great concerns on the negative impacts of DNG images, and some governments [14] have even amended laws to prevent the malicious sharing of fake media contents generated by machine learning software like DeepFake. In order to determine the authenticity of images and avoid the potential security issues, it is of importance to identify DNG images.

Identification of fake images is much related to the research field of information security, particularly in image forensics [15] and biometric anti-spoofing [16]. Due to the fact that fabricating a fake image would inevitably introduce some traces, the key to the identification is to analyze and extract features that represent the corresponding traces. For example, the quantization artifacts are used in JPEG image forensics [17, 18], the joint artifacts left by different image operations can be used to determine the operation chains [19, 20, 21], the splicing inconsistencies are exploited to locate tampered image regions [22, 23], and the displaying/imaging distortions are utilized in face spoofing detection [24, 25]. Recently, some works have been developed to identify fake images generated by deep networks. Some of them employ the visual artifacts [26] or saturation abnormalities [27] in DNG images, and some of them are based on deep learning paradigm [28, 29, 30, 31]. Although these works introduced some applicable approaches to identify DNG images, they did not analyze the common inherent traces left in generated images and thus failed to provide interpretable conclusions. Furthermore, they did not fully consider some challenging scenarios, *e.g.*, the image sources or the generative models are different in the training and testing phases, which need to be carefully coped with in practice. To this end, more efforts should be devoted to the identification of DNG images.

In this paper, we propose an effective and interpretable method to identify DNG images. Considering the difference between camera imaging and DNG image generation, we analyze the disparity between DNG images and real images

and observe that the DNG images are more distinguishable in residual domain of chrominance components. Based on the observations, we design a feature set to identify DNG images, which is composed by co-occurrence matrix of residual images in different color components. Furthermore, we address several challenging detection scenarios, including the cases when the training and testing data are mismatched and the GAN model is unaware. Extensive experimental results have shown the effectiveness of the proposed method: 1) When the training and testing data are matched, the proposed method achieves high accuracies. 2) When the training and testing images are generated by the same type of GAN but with different semantic content types, the proposed method outperforms existing methods; this situation has not been studied in previous works. 3) When the training and testing images are generated by different GANs, the proposed method also achieves superior performance, although the performance varies in different cases. 4) When the GAN model is unaware, the proposed method can obtain promising results with one-class classification. The contributions of this paper are summarized as follows.

- Considering that DNG image generation is different from camera imaging, we have conducted the analysis and demonstrated the disparities between DNG images and real images. By measuring the correlations between adjacent pixels in some color spaces, we have found that the statistical properties of DNG images and real images are different in the chrominance components of HSV and YCbCr color spaces, while the disparities are more distinct in the residual domain.
- We have proposed an effective feature set for DNG image identification. The feature set consists of co-occurrence matrices extracted from the residual images of several color components. The proposed feature set is of low dimension, and achieves good detection performance even with a small training set.
- We have evaluated the identification performance in different detection scenarios. The proposed method outperforms the existing methods when the image semantic types or the GAN models are mismatched in the training and testing phases. It is worthy to mention that, in the GAN model-unaware case (only real images are available for training), we can still achieve good results by performing one-class classification.

The rest of this paper is organized as follows. Section 2 introduces some related works. Section 3 presents the details of the proposed method. Section 4 reports and discusses the experimental results. Finally, the concluding remarks are drawn in Section 5.

2. Related Works

2.1. Image Generative Models and GANs

A generative model can be trained with some given data, and it aims to produce samples that follow the same distribution as the training data. In an

ideal case, by improving the model and increasing the amount and the quality of training data, the generative model is expected to eventually generate plausible samples similar to those coming from real world. Currently the most popular generative models are based on deep neural networks, including Generative Adversarial Networks (GANs) [1], Variational Autoencoders (VAEs) [2], and autoregressive models [3]. VAEs and autoregressive models usually produce images of poor quality compared to GANs, which would limit their applications. Hence, most of the state-of-the-art image generative models are built with GANs.

GAN was first proposed by Goodfellow *et al.* [1]. Basically, a GAN consists of two networks: a generator and a discriminator. The generator tries to generate synthetic samples as the one drawing from real data distribution, and the discriminator tries to correctly classify whether samples are coming from the generator or the real data. The training of GAN works as solving a two-player zero-sum game between the generator and the discriminator. While the discriminator notices some differences between the real distribution and the generated distribution, the generator adjusts its parameters to produce samples closer to the real distribution. And then, the discriminator tries to tell apart the two distributions again by adjusting its parameters. In an ideal case, the generator is expected to eventually reproduce the distribution of real data, and the discriminator fails to distinguish between generated samples and real samples. Up to now, many works [32, 33, 34, 35, 36, 4, 5, 6] have been proposed to improve the vanilla GAN. For example, Radford *et al.* [32] designed deep convolutional GAN (DCGAN), Arjovsky *et al.* [35] adopted Wasserstein distance in GAN to make the training more stable, Gulrajani *et al.* [36] made an improvement for Wasserstein GAN with gradient penalty (WGAN-GP). More recently, the quality and variation of generated images have been further improved by progressive growing of GANs (PROGAN) [4], large scale training with architectural changes and modified regularization scheme (BIGGAN) [5], designing style-based generator architecture (STYGAN) [6], and generating images by parts based on their conditional spatial coordinates (COCOGEN) [37].

In this paper, we consider DNG images generated from six popular GAN models, including DCGAN [32], WGAN-GP [36], PROGAN [4], STYGAN [6], BIGGAN [5], and COCOGEN [37]. While some of these GANs have already been considered in the recently-developed image database FaceForensic++ [38], we include some more different GAN models to evaluate the generalization capability of detection methods.

2.2. Fake Image Identification

Before the appearance of DNG images, many methods were developed to identify fake images subjected to editing and/or rebroadcasting. Most of those methods rely on specific traces regarding to the processing pipelines of fake images, such as the JPEG compression errors of a recompressed image [17] and the quality distortions of a spoofing face image [39, 40]. Hence, such targeted methods are not suitable for detecting DNG images. On the other hand, there are some general methods that are based on either statistical features extracted from image textures [25, 41] or deep neural networks equipped with a constrained

convolutional layer[42]. Such approaches may be used to detect DNG images, as long as retraining a detector with DNG image samples. Nevertheless, since these approaches do not take into account the traces left by the generation pipeline of DNG images, they cannot always achieve satisfactory performance.

With the growing interest in identifying DNG images, a few methods have been designed to differentiate between DNG images and real images. Some of these methods rely on feature engineering. McCloskey and Albright [27] found that the normalization of an image generator is different from a real camera, and proposed to detect GAN-generated images with saturation cues. The best result of the area under the ROC curve (AUC) obtained by this method was just around 0.70. Matern *et al.* [26] reported that some generated face images exhibited visual artifacts in eyes, teeth and facial contours, and utilized such artifacts for identification, obtaining the AUC around 0.85. Marra *et al.* [43] showed that each GAN leaves a specific fingerprint in the images generated by itself. Some low-level facial artifacts can be used to expose the images/videos generated by DeepFake, such as the absence of eye blinking [44], the errors of estimated 3-D head poses [45], the inconsistencies in facial parts locations [46], and the correlations of facial expressions and movements [47]. However, such facial abnormalities can be only applicable to generated face images/videos, and thus they are not suitable for detecting DNG images with other types of semantic contents.

Other methods resort to deep learning (DL). Marra *et al.* [28] tested the performance of several DL architectures for detecting images produced by image-to-image translation [9]. Mo *et al.* [29] and Dang *et al.* [30] designed customized convolutional neural networks (CNN) to identify fake face images generated by GANs. Afchar *et al.* [48] employed networks focusing on the mesoscopic properties of images to perform detection. In addition to directly feeding the images into CNNs, some works tried to improve the detection performance by incorporating specific domain knowledge. For example, training the network with co-occurrence matrices extracted from an image in the pixel domain of RGB space [49], prompting the network to learning the affine face warping artifacts [50] or up-sampling artifacts [51].

A big potential practical problem for identifying DNG images is the mismatch between the training and testing data. In order to mitigate the performance degradation of a trained detector in such case, researchers have tried to utilize more advanced learning mechanisms, *e.g.*, weakly-supervised transfer learning [52], incremental learning [53] and two-step pairwise learning [31]. However, these approaches still need to appropriately collect enough DNG images for the training procedure.

3. DNG Image Identification

In this section, we first analyze some possible artifacts of DNG images and investigate the disparities between DNG images and real images in some color spaces. Then, we construct a feature set to capture the artifacts of DNG images

so as to identify them. Finally, we discuss several detection scenarios and the corresponding detection strategies.

3.1. Analysis from the Perspective of Color

3.1.1. The generation pipeline of DNG images

In order to distinguish DNG images from real images, it needs to inspect the artifacts left by GANs during creating images. Typically, the generator of a GAN takes a random latent vector as input and employs several convolution/deconvolution layers to gradually expand the spatial size of the random vector. In the last layer of the generator, several feature maps are transformed into a tensor with three channels, where the three channels represent the R, G, and B components of the generated image, respectively. During this procedure, the convolution operations would introduce some inherent properties into the DNG images. By contrast, a real image is captured from real scene by camera, where the color components are decomposed and digitalized from real world, meaning that the real pixels should be inherently correlated in a different way. Since the generation of DNG images is quite different from camera imaging, there should be certain disparities between DNG images and real images from the perspective of color. Therefore, it is reasonable to assume that some properties among the color components of DNG images are different from real ones. In the following analysis, we will show some experimental evidences to support this assumption.

3.1.2. Discernibility of color component

As GANs usually generate images in RGB space, they tend to follow the properties of real images in RGB space, while paying less attention to the properties in other color spaces. In this way, although the differences between DNG images and real ones are inapparent in RGB color space, they may be more obvious in other color spaces. Therefore, we consider analyzing DNG images in three different color spaces, *i.e.*, RGB, HSV, and YCbCr, and try to use a metric to examine which color component is more discernible for identifying DNG images. Firstly, we construct the discernibility metric as follows.

- a) For the i -th image \mathbf{I} in a dataset, we calculate the correlation coefficient between the adjacent pixels in each of its color component \mathbf{I}^c ($c \in \{\text{R, G, B, H, S, V, Y, Cb, Cr}\}$), which can be formulated as:

$$r_i^c = \frac{\sum_{j=1}^m \sum_{k=1}^{n-1} (\mathbf{I}_{j,k}^c - \bar{\mathbf{I}}^c) (\mathbf{I}_{j,k+1}^c - \bar{\mathbf{I}}^c)}{\sqrt{\sum_{j=1}^m \sum_{k=1}^{n-1} (\mathbf{I}_{j,k}^c - \bar{\mathbf{I}}^c)^2 \sum_{j=1}^m \sum_{k=1}^{n-1} (\mathbf{I}_{j,k+1}^c - \bar{\mathbf{I}}^c)^2}}, \quad (1)$$

where $\bar{\mathbf{I}}^c$ is the mean value of \mathbf{I}^c , and m and n are the height and width of the image. In this way, the value of r_i^c represents the relevance of the adjacent pixel values. The larger the r_i^c , the higher correlation between the adjacent pixel values in \mathbf{I}^c . Please note that we only employ the correlation analysis on horizontal adjacent pixels for simplifying the analysis, and

similar results can be obtained by considering adjacent pixels in other directions.

- b) For a set of DNG images, we calculate r_i^c for each image and construct the histogram of r_i^c as $\mathbb{H}_{\text{DNG}}^c$. Similarly, for a set of real images, we construct the histogram $\mathbb{H}_{\text{Real}}^c$. Then, we measure the similarity between the two histograms with Chi-square distance

$$d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c) = \frac{1}{2} \sum_x \frac{(\mathbb{H}_{\text{DNG}}^c(x) - \mathbb{H}_{\text{Real}}^c(x))^2}{\mathbb{H}_{\text{DNG}}^c(x) + \mathbb{H}_{\text{Real}}^c(x)}, \quad (2)$$

where x is the bin index. $d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c)$ can serve as the *discernibility metric*. The larger the discernibility metric, the disparities between DNG images and real images are more significant. In other words, the color component c is more powerful for distinguishing between DNG images and real images with a larger value of $d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c)$.

To evaluate the effectiveness of the discernibility metric, we perform some analytical experiments. We use the STYGAN [6] model trained with the FFHQ dataset [6] as an example of image generative model. With the trained model, we generate a large amount of DNG images. We randomly select 10,000 DNG images and 10,000 real images from the generated dataset and the FFHQ dataset, respectively, and then compute the r_i^c ($c \in \{\text{R, G, B, H, S, V, Y, Cb, Cr}\}$) as introduced above. Then, we construct the histograms $\mathbb{H}_{\text{DNG}}^c$ and $\mathbb{H}_{\text{Real}}^c$. The histograms for different color components are shown in Fig. 1. It can be observed that the non-overlapping regions of $\mathbb{H}_{\text{DNG}}^c$ and $\mathbb{H}_{\text{Real}}^c$ are quite small for the R, G, and B channels, meaning that the disparities between DNG images and real images are not apparent in the RGB color space. This can be explained by the DNG generation process where R, G, B channels are directly generated to follow the distributions of real images. Via converting the images from RGB space into HSV/YCbCr spaces, the R, G, and B channels are merged by linear or non-linear combinations to form the other color channels. In this way, the disparities between DNG images and real images will be amplified, and thus the non-overlapping regions of $\mathbb{H}_{\text{DNG}}^c$ and $\mathbb{H}_{\text{Real}}^c$ become larger, especially for the four chrominance components, *i.e.*, H, S, Cb, and Cr. Such phenomena can be supported by the discernibility metrics, *i.e.*, $d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c)$. As shown in the sub-captions of the corresponding sub-figures in Fig. 1, we observe that the values of $d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c)$ for the four chrominance components are all exceed 0.011, which are larger than those for the luminance components (*i.e.*, V and Y) and R, G, and B components (all being less than 0.005). These results indicate that the chrominance components are more discernible than the others.

3.1.3. Discernibility analysis in first-order differential residual domain

In the above analysis, we directly derive the observations from the image spatial domain. However, the contents of DNG images and real images are not always visually distinguishable, meaning that the contents in DNG images and real images are quite similar. This would have negative impacts on analyzing

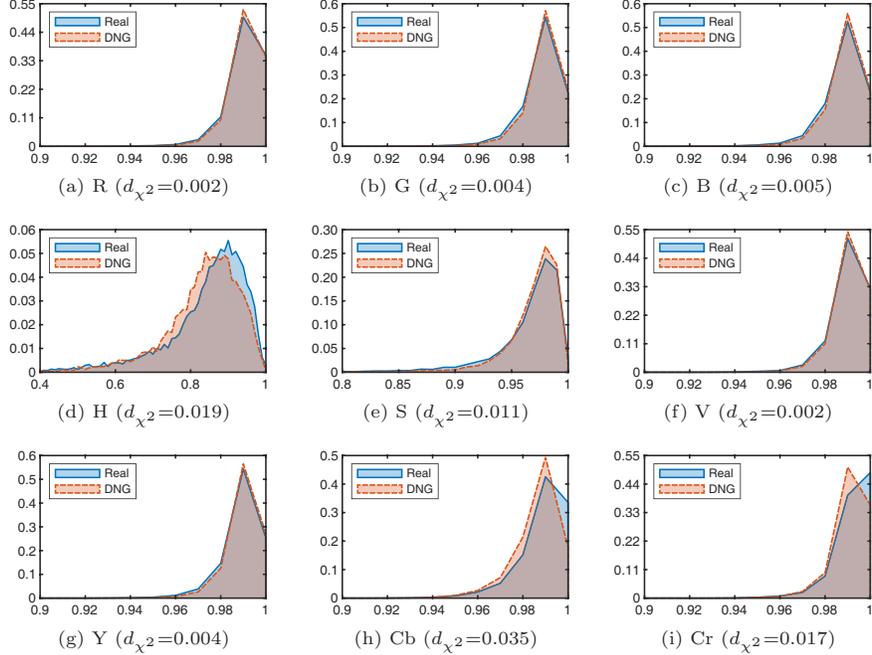


Figure 1: The histograms $\mathbb{H}_{\text{DNG}}^c$ (red) and $\mathbb{H}_{\text{Real}}^c$ (blue) for different color components. The values of $d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c)$ are included in the sub-captions.

their disparities. Therefore, it is more reasonable to suppress image contents with high-pass filtering and then investigate the disparities in image residuals. In fact, extracting features from high-pass filtered image residuals have been successfully used in some imperceptible pattern recognition applications, such as image steganalysis [54] and image forensics [41].

We apply a first-order differential operator as an example to obtain image residuals, namely,

$$\mathbf{R}_{j,k}^c = \mathbf{I}_{j,k}^c - \mathbf{I}_{j,k+1}^c, \quad c \in \{\text{R, G, B, H, S, V, Y, Cb, Cr}\}. \quad (3)$$

where \mathbf{I}^c is the c -th component of image \mathbf{I} and \mathbf{R}^c is the corresponding residual. We only consider horizontal difference here, and similar results can be obtained by considering vertical difference. Having the image residuals in residual domain, we replace \mathbf{I}^c and $\bar{\mathbf{I}}^c$ in Equation (1) with \mathbf{R}^c and $\bar{\mathbf{R}}^c$, respectively, and then conduct the discernibility analysis with the same image datasets as described in Section 3.1.2. The histograms $\mathbb{H}_{\text{DNG}}^c$ and $\mathbb{H}_{\text{Real}}^c$ for different color components in residual domain are shown in Fig. 2. From this figure, it is observed that the non-overlapping regions of $\mathbb{H}_{\text{DNG}}^c$ and $\mathbb{H}_{\text{Real}}^c$ become much more noticeable, implying that the DNG image and real image are indeed more distinguishable in the residual domain. This phenomenon can also be explained by the fact that DNG images are undergone a different processing pipeline from real

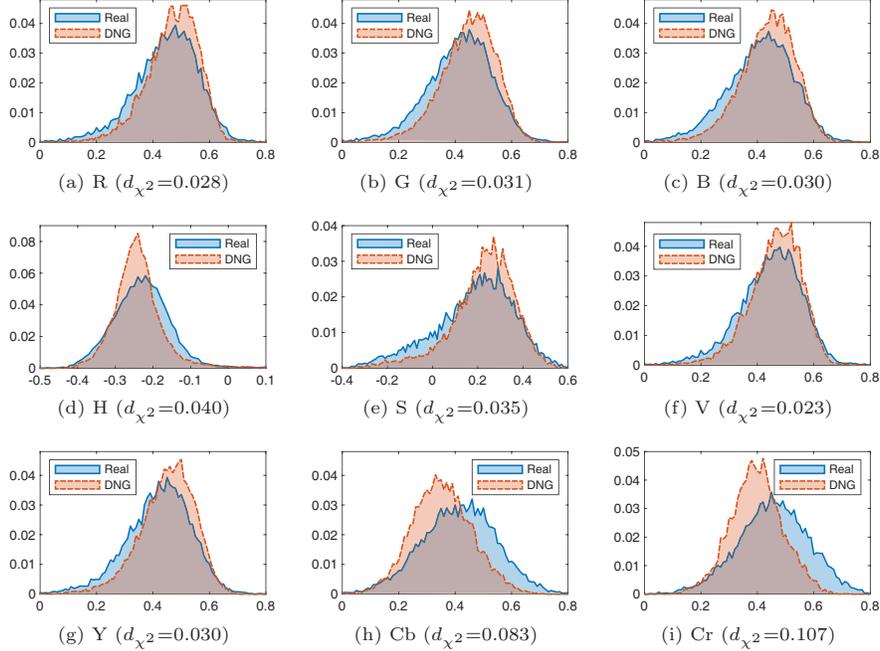


Figure 2: The histograms $\mathbb{H}_{\text{DNG}}^c$ (red) and $\mathbb{H}_{\text{Real}}^c$ (blue) for different color components in the residual domain. The values of $d_{\chi^2}(\mathbb{H}_{\text{DNG}}^c, \mathbb{H}_{\text{Real}}^c)$ are included in the sub-captions.

images so that their statistics are different in the residual domain. Moreover, the observations obtained in spatial domain still hold in residual domain: 1) the non-overlapping regions of $\mathbb{H}_{\text{DNG}}^c$ and $\mathbb{H}_{\text{Real}}^c$ for the H, S, Cb, Cr components are larger than those for R, G, B, V, Y components, and 2) the discernibility metrics for the four chrominance components are also greater than those for the other components. We have also conducted the same analysis on images generated by other types of GANs and obtained consistent results, indicating that some statistical features extracted from the chrominance components in residual domain would be beneficial for identifying DNG images.

3.2. Extracting Features from Color Components

Based on the previous analysis, we decide to extract features from the chrominance components in residual domain of an image. The overall framework of the proposed method is illustrated in Fig. 3. For a given image, we first compute the features from color components and then concatenate them into a feature vector, and finally train a classifier to predict whether the image is real or is generated by deep networks. In the feature extraction stage, we first calculate the image residuals with the first-order differential operators and pre-process them with truncation. Then, the co-occurrence matrix [55], which

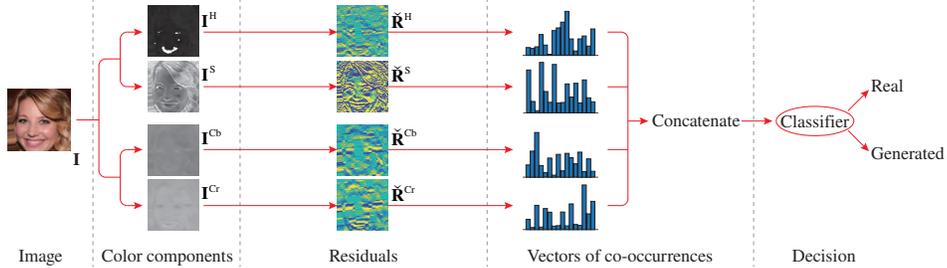


Figure 3: The overall framework of the proposed method.

has been widely used in image textural analysis as a kind of feature descriptor, is extracted from the image residuals and merged to form the feature set.

3.2.1. Truncating image residuals

Once the image residuals of H, S, Cb, Cr components have been obtained with the first-order differential operators, we need to pre-process the residuals before computing co-occurrence matrix. The reason is that there are too many distinct element values in the residuals, and it would result in a co-occurrence matrix with huge dimension if the raw residual data is directly used. In order to reduce the number of distinct values, the residual images \mathbf{R}^c ($c \in \{H, S, Cb, Cr\}$) are truncated as follow:

$$\check{\mathbf{R}}^c(x, y) = \begin{cases} \tau, & \mathbf{R}^c(x, y) \geq \tau, \\ \mathbf{R}^c(x, y), & -\tau < \mathbf{R}^c(x, y) < \tau, \\ -\tau, & \mathbf{R}^c(x, y) \leq -\tau, \end{cases} \quad (4)$$

where (x, y) is the position index of an element within the residual image and $\tau > 0$ is the truncation threshold. After truncation, the resulting residual images $\check{\mathbf{R}}^c$ only contain integer values within the range of $[-\tau, \tau]$. They are then used to compute the co-occurrence matrices.

3.2.2. Extracting co-occurrence features

In total, we have four co-occurrence matrices, which are calculated from $\check{\mathbf{R}}^H$, $\check{\mathbf{R}}^S$, $\check{\mathbf{R}}^{Cb}$, and $\check{\mathbf{R}}^{Cr}$, respectively. Typically, the co-occurrence matrix of a 2-D array \mathbf{V} is computed by

$$\mathbf{C}(\theta_1, \theta_2, \dots, \theta_d) = \frac{1}{n} \sum_{x, y} \mathbb{1} \left(\mathbf{V}(x, y) = \theta_1, \right. \\ \left. \mathbf{V}(x + \Delta x, y + \Delta y) = \theta_2, \dots, \right. \\ \left. \mathbf{V}(x + (d-1)\Delta x, y + (d-1)\Delta y) = \theta_d \right), \quad (5)$$

where $\mathbb{1}(\cdot)$ is an indicator function, $(\theta_1, \theta_2, \dots, \theta_d)$ is the index of co-occurrence matrix, d is the order of co-occurrence matrix, n is the normalization factor,

and Δx , Δy are the offsets for two neighboring elements. The dimensionality of each the co-occurrence matrix for $\tilde{\mathbf{R}}^c$ is $(2\tau + 1)^d$.

Since the co-occurrence matrix is symmetric, we can decrease the feature dimension by combining the two bins, $\mathbf{C}(\theta_1, \theta_2, \dots, \theta_d)$ and $\mathbf{C}(\theta_d, \theta_{d-1}, \dots, \theta_1)$, into one bin. After combination, the dimensionality of co-occurrence matrix is substantially decreased: the co-occurrence matrices for $\tilde{\mathbf{R}}^c$ ($c \in \{\text{H}, \text{S}, \text{Cb}, \text{Cr}\}$) have only $((2\tau + 1)^d + (2\tau + 1)^{d-1})/2$ bins. Please note that the reduction of feature dimension is motivated by the natural symmetry property of images, thus it would not significantly decrease the detection performance. In fact, the resulting relatively low dimensional features will speed up the training of classifier.

3.2.3. Practical implementation

In our practical implementation, we use two first-order differential operators, one in horizontal direction and one in vertical direction, to obtain the image residuals. The residuals are then processed as described above, where the truncation threshold is set as $\tau = 2$, and the order of co-occurrence matrix is set as $d = 3$. We choose the offsets as $(\Delta x, \Delta y) \in \{(0, 1), (1, 0)\}$, and thus we have 4 co-occurrence matrices (2 residuals \times 2 offsets) for each color component. Finally, we take the element-wise sum of the 4 co-occurrence matrices as the features. In total, a 300-D feature set is obtained, in which the feature dimension for each of $\tilde{\mathbf{R}}^{\text{H}}$, $\tilde{\mathbf{R}}^{\text{S}}$, $\tilde{\mathbf{R}}^{\text{Cb}}$, and $\tilde{\mathbf{R}}^{\text{Cr}}$ is $(5^3 + 5^2)/2 = 75$.

3.3. Detection Scenarios and Strategies

In practice, there are many types of generative models, and they can be trained with different real data for producing images with diverse contents. As a result, DNG images generated by different models that are trained with different datasets may not exhibit the same characteristic, leading to difficulties in distinguishing them from real images. Based on the information that is available in the training phase and the source of testing data, we divide the detection scenarios into three cases, and discuss the corresponding detection strategies as follows.

3.3.1. Matched training-testing data

In this case, the training and testing DNG images are generated by the same model trained with the same real image dataset. This is the simplest case and was considered by default in some existing works, *e.g.*, [28, 29, 27, 30]. To perform the detection, the investigator just need to train a binary classifier with real images and DNG images, and uses the trained classifier to predict the class labels for the testing images.

3.3.2. Mismatched training-testing data

In this case, the training and testing DNG images are from different sources. For example, they are generated by the same type of GAN but with different datasets or even different image semantic content types, or they have the same

type of image semantic contents but are generated with different GANs. Binary classification is also used in this case. Due to the mismatched sources, this case can be exploited to evaluate the generalization ability of a detection method. The better results a method achieves in this case, the better the method is expected to perform in real applications. Some existing works (*e.g.*, [52, 53, 31]) have been aware of this problem and tried to improve the performance with advanced learning mechanisms, but they have not intensively studied the performance for such mismatch cases.

3.3.3. Model-unaware case

It is not rare that the investigator does not have any knowledge about the generative model, and thus has no DNG image samples when building a classifier. This is the most challenging case presented to the investigator. To our best knowledge, such situation has never been considered in the existing works on DNG image identification. To cope with this case, a possible way is to train a one-class classifier with only real images and use the classifier to detect whether the testing image is real or not. If a testing image is not predicted as real, then it will be treated as a generated one.

4. Experiments

In this section, we will evaluate the performance of the proposed method. We first introduce the common experimental settings, and then present the experimental results for the three detection scenarios described in Section 3.3.

4.1. Experimental Setups

4.1.1. Real image datasets

Five real images datasets with two types of semantic contents (*i.e.*, face and bedroom) and different resolutions were used in the experiments. Each real image dataset is denoted as \mathcal{R}_α , where the subscript α indicates image source and resolution. The details of the real image datasets are described as follows, and their basic information can be referred to the upper part of Table 1.

- *Low-Resolution (LR) Face* (\mathcal{R}_{F-LR} and \mathcal{R}_{FI-LR}): The \mathcal{R}_{F-LR} dataset consists of 200,000 celebrity face images that were randomly selected from the “Align&Cropped” PNG images in the CelebA dataset [56]. We first cropped a 138×138 facial region from each image to remove the background, and then resized the cropped region to 128×128 . The \mathcal{R}_{FI-LR} dataset contains 10,000 face images in the LFW dataset [57]. These images are with the original size of 250×250 . For each image, we first cropped a 150×150 facial region and then resize it to 128×128 . Please note that since \mathcal{R}_{FI-LR} contains much fewer images than others, the images in this dataset (and the DNG images related to this datasets) are only used for testing in the data-mismatch case (Section 4.3.1).

Table 1: The real and generated datasets used in the experiments.

Dataset	Category	Content	Resolution	Quantity	Note
\mathcal{R}_{F-LR}	Real	Face	128×128	200,000	Selected from CelebA
\mathcal{R}_{F1-LR}	Real	Face	128×128	10,000	Selected from LFW
\mathcal{R}_{B-LR}	Real	Bedroom	128×128	200,000	Selected from LSUN bedroom
\mathcal{R}_{F-HR}	Real	Face	1024×1024	100,000	Combination of CelebA-HQ and FFHQ
\mathcal{R}_{B-HR}	Real	Bedroom	256×256	100,000	Selected from LSUN bedroom
\mathcal{G}_{F-LR}^* [†]	Generated	Face	128×128	200,000	GANs trained with CelebA
\mathcal{G}_{F1-LR}^*	Generated	Face	128×128	10,000	GANs trained with LFW
\mathcal{G}_{B-LR}^*	Generated	Bedroom	128×128	200,000	GANs trained with LSUN bedroom
\mathcal{G}_{F-HR}^*	Generated	Face	1024×1024	100,000	GANs trained with CelebA-HQ or FFHQ
\mathcal{G}_{B-HR}^*	Generated	Bedroom	256×256	100,000	GANs trained with bedroom images

[†] The asterisk * denotes the type of GAN. $*$ \in {DCGAN, WGAN-GP} for LR datasets, and $*$ \in {PROGAN, STYGAN, BIGGAN, COCOGAN} for HR bedroom datasets, and $*$ \in {PROGAN, STYGAN} for HR face datasets.

- *Low-Resolution Bedroom* (\mathcal{R}_{B-LR}): This dataset consists of 200,000 bedroom images that were randomly selected from the LSUN Bedroom dataset [58]. We first cropped the central region with the size of 256×256 from each image, and then resized the cropped region into 128×128 .
- *High-Resolution (HR) Face* (\mathcal{R}_{F-HR}): This dataset consists of 100,000 face images, including 30,000 images in the CelebA-HQ dataset [4] and 70,000 images in the FFHQ dataset [6]. These images are with the size of 1024×1024 , and they are resized to 256×256 in experiments so as to reduce the computational time (especially for deep learning based methods). Although the images have been pre-processed with down-scaling, the detection decision obtained by the classifier can also be applied to the original large images.
- *High-Resolution Bedroom* (\mathcal{R}_{B-HR}): This dataset consists of 100,000 bedroom images, which were randomly selected from the LSUN Bedroom dataset [58]. We cropped the central 256×256 region from each image.

4.1.2. DNG image datasets

In the experiments, we adopted DNG images generated by six types of GANs, including DCGAN [32], WGAN-GP [36], PROGAN [4], STYGAN [6], BIGGAN [5], and COCOGAN [37]. For DCGAN and WGAN-GP, we used the implementation available online¹ and trained GANs with the LR real image datasets (*i.e.*, \mathcal{R}_{F-LR} , \mathcal{R}_{F1-LR} , and \mathcal{R}_{B-LR}). Hence, the trained GANs produce LR DNG images with the size of 128×128 . For PROGAN, STYGAN, BIGGAN and COCOGAN, we used the corresponding trained models² released by the authors to generate

¹ Available at https://www.github.com/igul222/improved_wgan_training.

² Available at <https://drive.google.com/open?id=0B4qLcYyJmizONHFULTdYc051X0U>, <http://stylegan.xyz/drive>, <https://tfhub.dev/deepmind/biggan-128/2>, and <https://drive.google.com/drive/folders/1r-BvW6cVMHKJw-0wMI6mUepMkboWwWqN>, respectively.

HR images (*i.e.*, 1024×1024 face images and/or 256×256 bedroom images). The 1024×1024 face images were resized to 256×256 in our experiments. For simplicity, in the following context we denote a DNG images dataset as \mathcal{G}_α^β , where α represents image source and resolution, and β represents the type of GAN model. For example, $\mathcal{G}_{F-LR}^{DCGAN}$ denotes the LR CelebA face images generated by a DCGAN model. The basic information of the DNG image datasets are summarized in the bottom part of Table 1. In total, more than 1,400,000 DNG images are involved in our experiments.

4.1.3. Comparative study

We have compared the proposed method with four hand-crafted feature based methods and five deep learning based methods, including general-purpose methods and targeted methods. The details are as follows.

- SRM [54] (General-purpose hand-crafted feature, 34671-D): It is originally designed for image steganalysis. The features are extracted from image residuals obtained by a series of high-pass filters.
- Sub-SRM [41] (General-purpose hand-crafted feature, 714-D): This is a refined subset of SRM. It achieved good performance in image forensics.
- CoALBP+LPQ [25] (General-purpose hand-crafted feature, 19968-D): This method is designed for face spoofing detection. The feature set is composed of Co-Occurrence of Adjacent Local Binary Patterns and Local Phase Quantization.
- Sat-Cues [27] (Targeted hand-crafted feature for DNG image detection, 24-D): It uses the saturation cues to classify the real images and images generated by GANs.
- VGG-16 [59] (General-purpose DL based method): A famous CNN for image classification.
- ResNet_v2-50 [60] (General-purpose DL based method): Another famous CNN for image classification, which uses identity connections to improve performance.
- Mo *et al.* [29] (Targeted DL based method for DNG image detection): A CNN for detecting the images generated by PROGAN.
- CGFace [30] (Targeted DL based method for DNG image detection): A customized CNN for computer-generated face detection.
- TS-CDNN [31] (Targeted DL based method for DNG image detection): A coupled deep neural network (CDNN) with a two-step learning approach for detecting GAN generated face images.

Please note that we accordingly modified the fully connected layers of the original CNN networks if necessary, so as to ensure that they are compatible with the image sizes.

Table 2: Classification results (%) for **matched** training and testing data.

Method	FPR	FNR	ACC
	Average (Best/Worst)	Average (Best/Worst)	Average (Best/Worst)
Sub-SRM [41]	0.11 (0.00 / 0.82)	0.05 (0.00 / 0.40)	99.92 (100.0 / 99.39)
SRM [54]	0.03 (0.00 / 0.19)	0.01 (0.00 / 0.07)	99.98 (100.0 / 99.87)
CoALBP+LPQ [25]	0.01 (0.00 / 0.08)	0.01 (0.00 / 0.03)	99.99 (100.0 / 99.95)
Sat-Cues [27]	30.66 (17.77 / 42.76)	25.53 (17.59 / 38.59)	71.91 (80.69 / 59.33)
VGG-16 [59]	0.28 (0.01 / 1.32)	0.55 (0.01 / 2.40)	99.58 (99.99 / 98.14)
ResNet_v2-50 [60]	0.56 (0.01 / 1.99)	0.74 (0.01 / 2.57)	99.35 (99.99 / 98.04)
Mo <i>et al.</i> [29]	0.51 (0.00 / 3.51)	0.56 (0.00 / 2.43)	99.46 (100.0 / 97.03)
CGFace [30]	4.75 (0.04 / 37.81)	4.67 (0.01 / 41.78)	95.29 (99.97 / 60.20)
TS-CDNN [31]	7.04 (0.04 / 22.32)	10.47 (0.20 / 67.30)	91.24 (99.88 / 63.91)
Proposed	0.25 (0.00 / 1.96)	0.27 (0.00 / 2.14)	99.74 (100.0 / 97.95)

4.1.4. Training and testing protocol

Unless otherwise specified, we followed the descriptions below for training and testing. In each experiment, 50,000 pairs of DNG images and real images (25% of the LR images and 50% of the HR images) were randomly selected for training, while the remaining images were used for testing. For the hand-crafted feature based methods, we trained the ensembles of LDA (linear discriminative analysis) learners [61] as classifiers. The parameters were set as defaults used in [61]. For the DL based methods, we trained the networks with 90% of the training data and kept 10% for validation. The Momentum optimizer was used in training the network of Mo *et al.* and VGG-16, while the Adam optimizer was used in training ResNet_v2-50, CGFace, and TS-CDNN; the learning rates were searched within a certain range and the optimal ones were respectively selected for the networks. The networks were trained 50 epochs with a batch size of 64 (for the network of Mo *et al.*, VGG-16, ResNet_v2-50, and CGFace) or 128 (for TS-CDNN). We saved the models every 1 epoch and chose the ones with the best validation accuracies as the final models. In the testing stage, we computed and reported the false positive rate (FPR, the probability that real images are identified as generated ones), false negative rate (FNR, the probability that generated images are identified as real ones), and the overall accuracy (ACC).

4.2. Performance on Matched Training-testing Data

In this subsection, we evaluate the performance for identifying DNG images in the case that the training and testing data are matched. For each detection method, we respectively trained 10 classifiers to classify the fake images in the 10 DNG image datasets (*i.e.*, $\mathcal{G}_{F-LR}^{DCGAN}$, $\mathcal{G}_{F-LR}^{WGAN-GP}$, $\mathcal{G}_{B-LR}^{DCGAN}$, $\mathcal{G}_{B-LR}^{WGAN-GP}$, $\mathcal{G}_{F-HR}^{ProGAN}$, $\mathcal{G}_{F-HR}^{StyGAN}$, $\mathcal{G}_{B-HR}^{ProGAN}$, $\mathcal{G}_{B-HR}^{StyGAN}$, $\mathcal{G}_{B-HR}^{BigGAN}$, and $\mathcal{G}_{B-HR}^{COCO}$) and their corresponding real images. The results are summarized in Table 2, where the average, the best, and the worst results of FPR, FNR, and ACC obtained by the 10 classifiers for each method are reported. From Table 2, it is observed that all methods achieve an average accuracy over 90% except the method based on saturation cues. These results indicate that the DNG images can be accu-

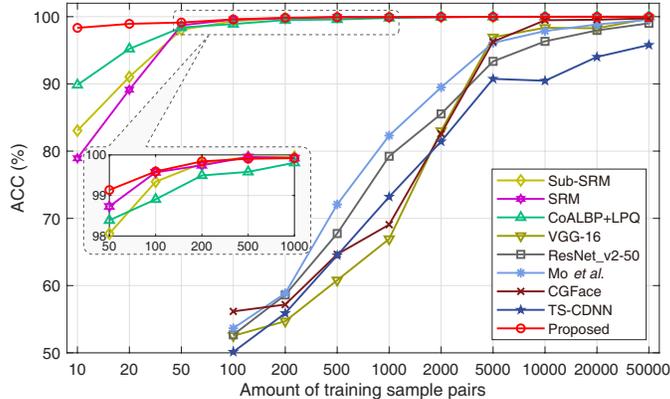


Figure 4: Classification accuracies for $\{\mathcal{R}_{F-LR}, \mathcal{G}_{F-LR}^{WGAN-GP}\}$ with different amounts of training sample pairs.

rately identified when a classifier is trained with corresponding data. Among all the methods, CoALBP+LPQ achieves the best average performance (99.99%), and the proposed method obtains competitive results (only 0.25% less). Please note that comparing with the high-dimensional CoALBP+LPQ feature, the proposed feature set is of much lower dimension, which can significantly reduce computational cost.

Another merit of the proposed feature is its advantages for small amount of training samples. We conducted an extra experiment to classify the images in \mathcal{R}_{F-LR} and $\mathcal{G}_{F-LR}^{WGAN-GP}$ by reducing training data. Fig. 4 shows the testing accuracies obtained with different amounts of training image pairs. It is observed that the proposed method achieves the testing accuracy of 98% even when using 10 pairs of training images, significantly outperforming the others in the same case. It means that the proposed method is less dependent on the amount of training data. When the number of training image pairs reaches 1,000, Sub-SRM, SRM, and CoALBP+LPQ can obtain almost the same performance with the proposed method, while the DL based methods only obtain accuracies that are less than 83%. According to Fig. 4, The DL based methods require more than 20,000 training image pairs to get competitive performance with the proposed method.

4.3. Performance on Mismatched Training-testing Data

In practice, the image source or the type of GAN model may not be fully known by the investigator. In this case, there is a mismatch between the training and testing data. In this subsection, we consider three sub-cases. In the first case, the training and testing images are from different sources and have the same type of semantic content, such as face or bedroom, while in the second case the image semantic content types are different in the training and testing phases. The last case is that the types of GAN models are mismatched.

Table 3: Classification results (%) for **mismatched image sources** (same semantic type).

	Method	FPR	FNR	ACC
<i>train:</i> $\{\mathcal{R}_{F-LR}, \mathcal{G}_{F-LR}^{DCGAN}\}$ <i>test:</i> $\{\mathcal{R}_{F1-LR}, \mathcal{G}_{F1-LR}^{DCGAN}\}$	Sub-SRM [41]	0.00	0.00	100.0
	SRM [54]	0.00	0.00	100.0
	CoALBP+LPQ [25]	0.00	0.00	100.0
	Sat-Cues [27]	42.68	10.65	73.34
	VGG-16 [59]	1.86	94.98	51.58
	ResNet_v2-50 [60]	0.00	100.0	50.00
	Mo <i>et al.</i> [29]	0.00	0.00	100.0
	CGFace [30]	1.25	99.02	49.87
	TS-CDNN [31]	0.04	99.85	50.06
	Proposed	0.00	0.00	100.0
<i>train:</i> $\{\mathcal{R}_{F-LR}, \mathcal{G}_{F-LR}^{WGAN-GP}\}$ <i>test:</i> $\{\mathcal{R}_{F1-LR}, \mathcal{G}_{F1-LR}^{WGAN-GP}\}$	Sub-SRM [41]	0.01	0.00	100.0
	SRM [54]	0.05	0.00	99.98
	CoALBP+LPQ [25]	0.00	0.11	99.95
	Sat-Cues [27]	65.54	25.06	54.70
	VGG-16 [59]	1.87	0.45	98.84
	ResNet_v2-50 [60]	8.29	87.98	51.87
	Mo <i>et al.</i> [29]	52.73	19.59	63.84
	CGFace [30]	0.64	1.30	99.03
	TS-CDNN [31]	14.83	39.76	72.71
	Proposed	0.00	2.50	98.75

4.3.1. Mismatched image sources (same semantic type)

It is reasonable to assume that the architecture of GAN is known, but the real images used to train the model are not available. To perform the detection, the investigator first trains a GAN with the same network architecture by using an alternative dataset, and then use the trained GAN to produce DNG samples. With these samples, the investigator can train a binary classifier to detect DNG images. In this experiment, we trained the binary classifiers with the LR face images related to the CelebA dataset (*i.e.*, \mathcal{R}_{F-LR} , $\mathcal{G}_{F-LR}^{DCGAN}$, and $\mathcal{G}_{F-LR}^{WGAN-GP}$), and then tested them with the LR face images related to the LFW dataset (*i.e.*, \mathcal{R}_{F1-LR} , $\mathcal{G}_{F1-LR}^{DCGAN}$, and $\mathcal{G}_{F1-LR}^{WGAN-GP}$). The detection results are shown in Table 3. It can be observed that all hand-crafted features except Sat-Cues achieve high accuracies in both cases, while the DL based methods perform poor in either of the two cases (VGG-16, Mo *et al.*, CGFace, and TS-CDNN), or both (ResNet_v2-50). The proposed method can perfectly detect the images generated by DCGAN even when they are from different sources, and also obtains very good result for detecting the images generated by WGAN-GP, though underperforms Sub-SRM by 1.25%.

4.3.2. Mismatched image sources (different semantic types)

To simulate the mismatch of image semantic contents, we first used the classifiers trained with face images to classify the bedroom images, and then used the classifiers trained with bedroom images to classify the face images. For example, a classifier trained with \mathcal{R}_{F-LR} and $\mathcal{G}_{F-LR}^{DCGAN}$ is used to classify the

Table 4: Classification results (%) for **mismatched image sources** (different semantic types).

Method	FPR	FNR	ACC
	Average (Best/Worst)	Average (Best/Worst)	Average (Best/Worst)
Sub-SRM [41]	23.72 (0.00 / 99.96)	17.54 (0.00 / 84.88)	79.37 (99.94 / 50.02)
SRM [54]	23.04 (0.00 / 99.58)	17.77 (0.00 / 99.00)	79.59 (100.0 / 50.21)
CoALBP+LPQ [25]	28.08 (0.00 / 98.77)	12.86 (0.00 / 88.28)	79.53 (99.99 / 49.74)
Sat-Cues [27]	42.72 (16.11 / 62.53)	45.54 (16.64 / 78.97)	55.87 (64.15 / 50.46)
VGG-16 [59]	19.81 (0.01 / 79.47)	61.26 (3.74 / 99.98)	59.46 (91.05 / 49.99)
ResNet_v2-50 [60]	4.30 (0.00 / 23.40)	95.19 (75.28 / 99.99)	50.26 (51.34 / 49.64)
Mo <i>et al.</i> [29]	17.45 (0.02 / 87.48)	50.01 (0.08 / 99.99)	66.27 (99.95 / 48.48)
CGFace [30]	16.77 (0.07 / 44.16)	63.87 (0.82 / 99.85)	59.68 (81.97 / 45.96)
TS-CDNN [31]	11.71 (0.01 / 28.37)	79.40 (49.92 / 99.97)	54.45 (72.13 / 50.02)
Proposed	28.33 (0.00 / 96.29)	9.45 (0.00 / 41.53)	81.11 (99.99 / 51.85)

images in \mathcal{R}_{B-LR} and $\mathcal{G}_{B-LR}^{DCGAN}$, and another classifier trained with \mathcal{R}_{B-LR} and $\mathcal{G}_{B-LR}^{DCGAN}$ is used to classify the images in \mathcal{R}_{F-LR} and $\mathcal{G}_{F-LR}^{DCGAN}$. The testing for LR images and HR images were conducted separately, and the DNG images in training and testing stages are produced by the same type of GAN, so there are totally 8 cases (Please note that the DNG images generated by BIGGAN and COCOGAN were not included in this experiment since they are all bedroom images). The average, best, and worst results are shown in Table 4. We observe that the DL based methods significantly underperform most of the feature based methods, implying that they over-fit on image contents and thus have poor generalization performance for images with different semantic types. The proposed method obtains lower FNR than others and achieves the best average performance (81.11%), outperforming the second place (*i.e.*, SRM) by 1.5%. On the other hand, it is observed that the detection performance varies greatly in different cases. For example, the proposed method achieves 99% accuracy on the best case (testing \mathcal{R}_{B-LR} and $\mathcal{G}_{B-LR}^{DCGAN}$ with the classifier trained on \mathcal{R}_{F-LR} and $\mathcal{G}_{F-LR}^{DCGAN}$), while it behaves close to random guessing in the worst case (testing \mathcal{R}_{F-HR} and $\mathcal{G}_{F-HR}^{ProGAN}$ with the classifier trained on \mathcal{R}_{B-HR} and $\mathcal{G}_{B-HR}^{ProGAN}$).

4.3.3. Mismatched GAN models

To simulate the mismatch of GAN models, we used the classifiers trained with the images generated by a certain GAN to classify the images generated by another GAN. The experiments for LR images and HR images were also conducted separately, and thus there are totally 18 cases (2 for LR face images, 2 for LR bedroom images, 2 for HR face images, and 12 for HR bedroom images). The experimental results are shown in Table 5. On average, the proposed method performs the best (91.87%) among all the methods, since it obtains good accuracies for testing the HR bedroom images even when the GAN models are mismatched. However, it does not always have good performance. In the worst case, namely, using the classifier trained with \mathcal{R}_{F-LR} and $\mathcal{G}_{F-LR}^{DCGAN}$ to test the images in \mathcal{R}_{F-LR} and $\mathcal{G}_{F-LR}^{WGAN-GP}$, all the DNG images are mistakenly classified, resulting in the accuracy of 50%. The reason may be that DCGAN is

Table 5: Classification results (%) for **mismatched GAN models** in training and testing.

Method	FPR	FNR	ACC
	Average (Best/Worst)	Average (Best/Worst)	Average (Best/Worst)
Sub-SRM [41]	0.06 (0.00 / 0.82)	39.45 (0.00 / 100.0)	80.25 (100.0 / 50.00)
SRM [54]	0.02 (0.00 / 0.19)	35.11 (0.00 / 100.0)	82.44 (100.0 / 50.00)
CoALBP+LPQ [25]	0.01 (0.00 / 0.08)	22.05 (0.00 / 100.0)	88.97 (100.0 / 49.99)
Sat-Cues [27]	26.79 (17.77 / 42.76)	40.18 (17.19 / 63.48)	66.52 (78.43 / 48.82)
VGG-16 [59]	0.34 (0.01 / 1.32)	77.09 (0.35 / 99.94)	61.29 (99.77 / 49.92)
ResNet_v2-50 [60]	0.54 (0.01 / 1.99)	82.09 (12.82 / 100.0)	58.68 (93.56 / 49.64)
Mo <i>et al.</i> [29]	0.32 (0.00 / 3.51)	58.04 (0.00 / 99.99)	70.82 (99.95 / 48.39)
CGFace [30]	6.87 (0.04 / 37.81)	71.18 (0.20 / 99.91)	60.98 (99.84 / 48.54)
TS-CDNN [31]	6.96 (0.04 / 22.32)	84.19 (55.68 / 99.96)	54.42 (68.16 / 46.81)
Proposed	0.14 (0.00 / 1.96)	16.12 (0.00 / 100.0)	91.87 (100.0 / 50.00)

simpler than WGAN-GP, and thus the images generated by DCGAN are more distinguishable than those generated by WGAN-GP, leading to that the trained classifier is too simple to identify the images generated by WGAN-GP. This indicates that the DNG samples must be carefully selected during training a binary classifier, otherwise the trained detector would produce wrong decisions in some mismatch cases.

4.4. Performance on Model-unaware One-class Classification

As mentioned above, the performance of a binary classifier for identifying DNG images may vary when the training and testing data are mismatch, and thus it is important to train the binary classifier with compatible DNG images. However, in practice the GAN model is sometimes unknown and there are no suitable DNG images available for training, thus it is difficult to build a feasible binary classifier. To perform the detection, an alternative way is to build a one-class classifier, which fits a model via learning the properties of real images and regard the DNG images as outliers. In this subsection, we evaluate the performance of the proposed method in one-class classification. Since it is not a trivial task to adapt a deep network for one-class classification, we only considered the feature based methods in this experiment (Sat-Cues was also excluded for its poor performance). We employed the support vector machine (SVM) with Gaussian kernel implemented in LIBSVM [62] as one-class classifier. The parameter γ of the Gaussian kernel was determined via a grid search. The parameter ν , which controls the upper bound of training error (*i.e.*, the probability of that regarding the training real images as outliers), was set as 0.10 and 0.05, respectively. In this experiment, the real face images and real bedroom images are combined together, and thus there are a LR real image dataset $\mathcal{R}_{F-LR} + \mathcal{R}_{B-LR}$ and a HR real image dataset $\mathcal{R}_{F-HR} + \mathcal{R}_{B-HR}$. Hence, we trained four classifiers for each method, corresponding to the two values of ν and LR/HR images, respectively. In the training phase, 50,000 real images were randomly selected to train the one-class SVMs. In the testing phase, the remaining real images and all DNG images were fed to the trained models to calculate the accuracies.

Table 6: Detection accuracies (%) obtained by **one-class classification**.

v	Train dataset	Method	$\mathcal{R}_{\text{F-LR}}^+$ $\mathcal{R}_{\text{B-LR}}$	Best	Worst	Average
0.10	$\mathcal{R}_{\text{F-LR}}$	Sub-SRM [41]	88.99	100.0	12.56	65.94
		SRM [54]	89.81	100.0	2.15	53.32
	$\mathcal{R}_{\text{B-LR}}$	CoALBP+LPQ [25]	89.64	100.0	13.49	67.80
		Proposed	89.90	100.0	99.85	99.94
0.05	$\mathcal{R}_{\text{F-LR}}$	Sub-SRM [41]	93.76	100.0	8.00	62.82
		SRM [54]	94.93	100.0	0.55	51.51
	$\mathcal{R}_{\text{B-LR}}$	CoALBP+LPQ [25]	94.66	100.0	4.35	59.45
		Proposed	94.91	100.0	99.19	99.69
v	Train dataset	Method	$\mathcal{R}_{\text{F-HR}}^+$ $\mathcal{R}_{\text{B-HR}}$	Best	Worst	Average
0.10	$\mathcal{R}_{\text{F-HR}}$	Sub-SRM [41]	89.65	99.95	4.44	46.38
		SRM [54]	90.15	77.44	5.67	30.71
	$\mathcal{R}_{\text{B-HR}}$	CoALBP+LPQ [25]	90.04	67.62	1.33	27.04
		Proposed	89.86	98.09	4.89	49.82
0.05	$\mathcal{R}_{\text{F-HR}}$	Sub-SRM [41]	94.29	99.81	1.92	41.26
		SRM [54]	95.03	64.14	2.56	22.66
	$\mathcal{R}_{\text{B-HR}}$	CoALBP+LPQ [25]	95.20	50.51	0.45	17.60
		Proposed	94.78	94.57	2.02	39.49

[†] $\mathcal{R}_{\text{F-LR}} + \mathcal{R}_{\text{B-LR}}$ denotes the combination of LR real image datasets $\mathcal{R}_{\text{F-LR}}$ and $\mathcal{R}_{\text{B-LR}}$, and $\mathcal{R}_{\text{F-HR}} + \mathcal{R}_{\text{B-HR}}$ denotes the combination of HR real image datasets $\mathcal{R}_{\text{F-HR}}$ and $\mathcal{R}_{\text{B-HR}}$.

The classification results are shown in Table 6, where the top half part contains the results for LR images and the bottom half part contains the results for HR images. In this table, we report the accuracies for detecting real images, as well as the average, best, and worst results for detecting DNG images. On the one hand, we observe that the detection performance for real images is related to the parameter v . Specifically, when v is 0.10 or 0.05, the testing errors for real images are correspondingly around 10% or 5%, respectively. It means that the trained classifier can adequately model the distribution of real images. On the other hand, it is observed that the proposed method can accurately detect the LR DNG images, outperforming the other methods by a large margin. For the HR DNG images, although our method obtains poor result in the worst case (for detecting the HR face images generated by PROGAN), it achieves relatively good performance on average. Based on these experimental results, it is promising to employ one-class classification with the proposed features to identify DNG images when only real images are available for training.

5. Conclusion

In this paper, we have investigated some new issues on identifying deep network generated images. We have analyzed and observed that the disparities between DNG images and real images are apparent in the residual domain of the

chrominance color components. Based on the observations, a feature set based on color statistical features is proposed. The feature set is compact and effective. We have evaluated the performance in some situations with different assumption on the availability of training data. The experimental results show that the proposed features equipped with a binary classifier can accurately differentiate between DNG images and real images when the training and testing data are matched, and outperforms existing methods when the image semantic types or the GAN models are mismatched in the training and testing phases. Moreover, in the model-unaware case, the proposed features can still effectively identify some types of DNG images with a one-class classifier.

In addition to the proposed detection method, this paper also provides some useful insights for the research community. Typically, the generative models generate images by imitating real images in RGB color space. Although the generated image may be visually appealing and be indistinguishable from real images for human eyes, they can be easily detected by the proposed method. It means that many inherent properties of real images, such as the properties in different color components, have not been properly depicted by the existing generative models. In order to further improve the quality of DNG images, more constrains should be considered in generative models.

In the future, we will improve the performance of our method to meet the requirements in practical detection situations. Furthermore, we will try to detect the modern image processing techniques that utilize deep networks based generative models as backbones, for example, image inpainting with GANs.

Acknowledgements

This work was supported in part by NSFC (61802262, 61872244 and U19B2022), Guangdong Basic and Applied Basic Research Foundation (2019B151502001), Guangdong R&D Program in Key Areas (2019B010139003), and Shenzhen R&D Program (JCYJ20180305124325555 and GJHZ20180928155814437).

References

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Proc. Conf. Neural Information Processing Systems (NeurIPS), 2014, pp. 2672–2680.
- [2] D. P. Kingma, M. Welling, Auto-encoding variational Bayes, in: Proc. Int. Conf. Learning Representations (ICLR), 2014.
- [3] A. van den Oord, N. Kalchbrenner, K. Kavukcuoglu, Pixel recurrent neural networks, in: Proc. Int. Conf. Machine Learning (ICML), 2016, pp. 1747–1756.

- [4] T. Karras, T. Aila, S. Laine, J. Lehtinen, Progressive growing of GANs for improved quality, stability, and variation, in: Proc. Int. Conf. Learning Representations (ICLR), 2018.
- [5] A. Brock, J. Donahue, K. Simonyan, Large scale GAN training for high fidelity natural image synthesis, arXiv preprint arXiv:1809.11096 (2018).
- [6] T. Karras, S. Laine, T. Aila, A style-based generator architecture for generative adversarial networks, arXiv preprint arXiv:1812.04948 (2018).
- [7] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al., Photo-realistic single image super-resolution using a generative adversarial network, in: Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2017, pp. 4681–4690.
- [8] Z. Wang, B. Chen, H. Zhang, H. Liu, Variational probabilistic generative framework for single image super-resolution, Signal Processing 156 (2019) 92–105.
- [9] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Proc. IEEE Int. Conf. Computer Vision (ICCV), 2017, pp. 2223–2232.
- [10] M.-Y. Liu, T. Breuel, J. Kautz, Unsupervised image-to-image translation networks, in: Proc. Conf. Neural Information Processing Systems (NeurIPS), 2017, pp. 700–708.
- [11] S. Iizuka, E. Simo-Serra, H. Ishikawa, Globally and locally consistent image completion, ACM Trans. Graphics 36 (4) (2017) 107:1–107:14.
- [12] J. Snow, AI could set us back 100 years when it comes to how we consume news, available: <https://www.technologyreview.com/s/609358> (Nov. 2017).
- [13] W. Knight, The us military is funding an effort to catch deepfakes and other AI trickery, available: <https://www.technologyreview.com/s/611146> (Nov. 2018).
- [14] Virginia bans ‘deepfakes’ and ‘deepnudes’ pornography, available: <https://www.bbc.com/news/technology-48839758> (Jul. 2019).
- [15] P. Korus, Digital image integrity—a survey of protection and verification techniques, Digital Signal Processing 71 (2017) 1–26.
- [16] J. Galbally, S. Marcel, J. Fierrez, Biometric antispoofing methods: A survey in face recognition, IEEE Access 2 (2014) 1530–1552.
- [17] W. Luo, J. Huang, G. Qiu, JPEG error analysis and its applications to digital image forensics, IEEE Trans. Inf. Forensics Security 5 (3) (2010) 480–491.

- [18] T. Bianchi, A. Piva, Image forgery localization via block-grained analysis of JPEG artifacts, *IEEE Trans. Inf. Forensics Security* 7 (3) (2012) 1003–1017.
- [19] Z. Chen, Y. Zhao, R. Ni, Detection of operation chain: JPEG-resampling-JPEG, *Signal Processing: Image Communication* 57 (2017) 8–20.
- [20] Q. Zhang, W. Lu, T. Huang, S. Luo, Z. Xu, Y. Mao, On the robustness of JPEG post-compression to resampling factor estimation, *Signal Processing* (2019) 107371.
- [21] M. C. Stamm, X. Chu, K. R. Liu, Forensically determining the order of signal processing operations, in: *Proc. IEEE Int. Workshop Information Forensics and Security (WIFS)*, 2013, pp. 162–167.
- [22] P. Korus, J. Huang, Multi-scale fusion for improved localization of malicious tampering in digital images, *IEEE Trans. Image Process.* 25 (3) (2016) 1312–1326.
- [23] H. Li, W. Luo, X. Qiu, J. Huang, Image forgery localization via integrating tampering possibility maps, *IEEE Trans. Inf. Forensics Security* 12 (5) (2017) 1240–1252.
- [24] K. Patel, H. Han, A. K. Jain, Secure face unlock: Spoof detection on smartphones, *IEEE Trans. Inf. Forensics Security* 11 (10) (2016) 2268–2283.
- [25] Z. Boulkenafet, J. Komulainen, A. Hadid, Face spoofing detection using colour texture analysis, *IEEE Trans. Inf. Forensics Security* 11 (8) (2016) 1818–1830.
- [26] F. Matern, C. Riess, M. Stamminger, Exploiting visual artifacts to expose deepfakes and face manipulations, in: *Proc. IEEE Winter Applications of Computer Vision Workshops (WACVW)*, 2019, pp. 83–92.
- [27] S. McCloskey, M. Albright, Detecting GAN-generated imagery using saturation cues, in: *Proc. IEEE Int. Conf. Image Processing (ICIP)*, 2019, pp. 4584–4588.
- [28] F. Marra, D. Gagnaniello, D. Cozzolino, L. Verdoliva, Detection of GAN-generated fake images over social networks, in: *Proc. IEEE Conf. Multimedia Information Processing and Retrieval (MIPR)*, 2018, pp. 384–389.
- [29] H. Mo, B. Chen, W. Luo, Fake faces identification via convolutional neural network, in: *Proc. ACM Workshop Information Hiding and Multimedia Security*, 2018, pp. 43–47.
- [30] L. M. Dang, S. I. Hassan, S. Im, J. Lee, S. Lee, H. Moon, Deep learning based computer generated face identification using convolutional neural network, *Applied Sciences* 8 (12) (2018) 2610.

- [31] Y.-X. Zhuang, C.-C. Hsu, Detecting generated image based on a coupled network with two-step pairwise learning, in: Proc. IEEE Int. Conf. Image Processing (ICIP), 2019, pp. 3212–3216.
- [32] A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional generative adversarial networks, in: Proc. Int. Conf. Learning Representations (ICLR), 2016.
- [33] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, S. P. Smolley, Least squares generative adversarial networks, in: Proc. IEEE Int. Conf. Computer Vision (ICCV), 2017, pp. 2813–2821.
- [34] J. Zhao, M. Mathieu, Y. LeCun, Energy-based generative adversarial network, in: Proc. Int. Conf. Learning Representations (ICLR), 2017.
- [35] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein GAN, arXiv preprint arXiv:1701.07875 (2017).
- [36] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, A. C. Courville, Improved training of Wasserstein GANs, in: Proc. Conf. Neural Information Processing Systems (NeurIPS), 2017, pp. 5769–5779.
- [37] C. H. Lin, C.-C. Chang, Y.-S. Chen, D.-C. Juan, W. Wei, H.-T. Chen, COCO-GAN: Generation by parts via conditional coordinating, in: Proc. IEEE Int. Conf. Computer Vision (ICCV), 2019, pp. 4512–4521.
- [38] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, M. Nießner, Faceforensics++: Learning to detect manipulated facial images, in: Proc. IEEE Int. Conf. Computer Vision (ICCV), 2019, pp. 1–11.
- [39] J. Galbally, S. Marcel, Face anti-spoofing based on general image quality assessment, in: Proc. Int. Conf. Pattern Recognition (ICPR), 2014, pp. 1173–1178.
- [40] D. Wen, H. Han, A. K. Jain, Face spoof detection with image distortion analysis, IEEE Trans. Inf. Forensics Security 10 (4) (2015) 746–761.
- [41] H. Li, W. Luo, X. Qiu, J. Huang, Identification of various image operations using residual-based features, IEEE Trans. Circuits Syst. Video Technol. 28 (1) (2018) 31–45.
- [42] B. Bayar, M. C. Stamm, Constrained convolutional neural networks: A new approach towards general purpose image manipulation detection, IEEE Trans. Inf. Forensics Security 13 (11) (2018) 2691–2706.
- [43] F. Marra, D. Gragnaniello, L. Verdoliva, G. Poggi, Do GANs leave artificial fingerprints?, in: Proc. IEEE Conf. Multimedia Information Processing and Retrieval (MIPR), 2019, pp. 506–511.

- [44] Y. Li, M.-C. Chang, S. Lyu, In *ictu oculi*: Exposing AI created fake videos by detecting eye blinking, in: Proc. IEEE Int. Workshop Information Forensics and Security (WIFS), 2018, pp. 1–7.
- [45] X. Yang, Y. Li, S. Lyu, Exposing deep fakes using inconsistent head poses, in: Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP), 2019, pp. 8261–8265.
- [46] X. Yang, Y. Li, H. Qi, S. Lyu, Exposing GAN-synthesized faces using landmark locations, in: Proc. ACM Workshop Information Hiding and Multimedia Security, 2019, pp. 113–118.
- [47] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, H. Li, Protecting world leaders against deep fakes, in: Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) Workshops, 2019.
- [48] D. Afchar, V. Nozick, J. Yamagishi, I. Echizen, Mesonet: a compact facial video forgery detection network, in: Proc. IEEE Int. Workshop Information Forensics and Security (WIFS), 2018, pp. 1–7.
- [49] L. Nataraj, T. M. Mohammed, B. Manjunath, S. Chandrasekaran, A. Flenner, J. H. Bappy, A. K. Roy-Chowdhury, Detecting GAN generated fake images using co-occurrence matrices, *Electronic Imaging 2019* (5) (2019) 532–1–532–7.
- [50] Y. Li, S. Lyu, Exposing deepfake videos by detecting face warping artifacts, in: Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) Workshops, 2019.
- [51] X. Zhang, S. Karaman, S.-F. Chang, Detecting and simulating artifacts in GAN fake images, in: Proc. IEEE Int. Workshop Information Forensics and Security (WIFS), 2019.
- [52] D. Cozzolino, J. Thies, A. Rössler, C. Riess, M. Nießner, L. Verdoliva, Forensictransfer: Weakly-supervised domain adaptation for forgery detection, *arXiv preprint arXiv:1812.02510* (2018).
- [53] F. Marra, C. Saltori, G. Boato, L. Verdoliva, Incremental learning for the detection and classification of GAN-generated images, *arXiv preprint arXiv:1910.01568* (2019).
- [54] J. Fridrich, J. Kodovsky, Rich models for steganalysis of digital images, *IEEE Trans. Inf. Forensics Security* 7 (3) (2012) 868–882.
- [55] R. M. Haralick, K. Shanmugam, I. Dinstein, Textural features for image classification, *IEEE Trans. Syst., Man, Cybern. SMC-3* (6) (1973) 610–621.
- [56] Z. Liu, P. Luo, X. Wang, X. Tang, Deep learning face attributes in the wild, in: Proc. IEEE Int. Conf. Computer Vision (ICCV), 2015, pp. 3730–3738.

- [57] G. B. Huang, M. Ramesh, T. Berg, E. Learned-Miller, Labeled Faces in the Wild: A database for studying face recognition in unconstrained environments, Tech. Rep. 07-49, University of Massachusetts, Amherst (Oct. 2007).
- [58] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, J. Xiao, LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop, arXiv preprint arXiv:1506.03365 (2015).
- [59] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556 (2014).
- [60] K. He, X. Zhang, S. Ren, J. Sun, Identity mappings in deep residual networks, in: Proc. European Conf. on Computer Vision (ECCV), 2016, pp. 630–645.
- [61] J. Kodovsky, J. Fridrich, V. Holub, Ensemble classifiers for steganalysis of digital media, *IEEE Trans. Inf. Forensics Security* 7 (2) (2012) 432–444.
- [62] C.-C. Chang, C.-J. Lin, LIBSVM: A library for support vector machines, *ACM Trans. Intell. Syst. Technol.* 2 (2011) 27:1–27:27, software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.