# Detection of Deep Network Generated Images Using Disparities in Color Components

Haodong Li, Bin Li, Shunquan Tan, Jiwu Huang

#### **Abstract**

With the powerful deep network architectures, such as generative adversarial networks and variational autoencoders, large amounts of photorealistic images can be generated. The generated images, already fooling human eyes successfully, are not initially targeted for deceiving image authentication systems. However, research communities as well as public media show great concerns on whether these images would lead to serious security issues. In this paper, we address the problem of detecting deep network generated (DNG) images by analyzing the disparities in color components between real scene images and DNG images. Existing deep networks generate images in RGB color space and have no explicit constrains on color correlations; therefore, DNG images have more obvious differences from real images in other color spaces, such as HSV and YCbCr, especially in the chrominance components. Besides, the DNG images are different from the real ones when considering red, green, and blue components together. Based on these observations, we propose a feature set to capture color image statistics for detecting the DNG images. Moreover, three different detection scenarios in practice are considered and the corresponding detection strategies are designed. Extensive experiments have been conducted on face image datasets to evaluate the effectiveness of the proposed method. The experimental results show that the proposed method is able to distinguish the DNG images from real ones with high accuracies.

#### **Index Terms**

Image generative model, generative adversarial networks, fake image identification, image statistics.

#### I. Introduction

With the rapid development of image processing technology, one can easily create image forgeries without leaving visual artifacts. The spread of fake images may result in moral, ethical, and legal consequences. It is important to identify fake information in order to avoid potential security issues. Therefore, determining the authenticity of images has attracted increasing attention in many applications, such as image forensics [1], [2] and biometric anti-spoofing [3].

Fabricating a fake image usually includes editing and/or rebroadcasting. Such a process would inevitably introduce some artifacts, which can be used for identifying the authenticity. For example, the quantization artifacts are used in JPEG image forensics [4], [5], the splicing inconsistences are used to detect the locations of tampered regions [6], [7], and the displaying/imaging distortions are utilized in face spoofing detection [8], [9].

Generative models have been widely used in many applications, such as speech synthesis [10], image super-resolution [11], image translation [12], [13], and image inpainting [14]. In the arsenal of fake image detection, however, there are few methods for identifying the faking process that creates images with generative models. The reason may be due to the fact that traditional generative methods can only generate simple image textures and the image contents are far from realistic. Therefore, it was not difficult to differentiate generative images in the last decades. In recent years, the situation starts to invert. With the advancement of deep learning, tremendous progress has been made in image generative models. Deep network generated (DNG) images become more and more photorealistic, and it is no longer easy to identify them with human eyes, which would lead to serious security risks. In fact, both public media [15] and research communities [16] recently have shown great concerns on the negative impacts of DNG images. As the training of generative models become more stable [17], [18] and the quality of DNG images become more satisfactory [19], the DNG images may be maliciously used for seeking illegal benefits. For example, some generated scenes can be used as materials to falsify images or videos that fabricate fake news; some generated faces can be posted on social networks by frauds to counterfeit personal information. Therefore, it is of importance to identify the DNG images. We try to address this problem in this paper. Although various contents can be created in DNG images, we mainly focus on detecting the generated face images, because most of existing generative models show great success in generating faces, and the generated faces may result in many security-related issues.

In this paper, we have analyzed the differences between DNG images and real images, and observe some phenomenons useful for differentiation. Specifically, we have observed that the DNG image is more differentiable in chrominance components. Based on these studies, we have proposed a method to detect DNG images by using color features extracted from image high-frequency parts. The features are composed by co-occurrence matrix of image high-pass filtering residuals in different color components. Besides, according to the information available to a detector, we have divided the detection scenarios into three different cases, *i.e.*, sample-aware, model-aware, and model-unaware. Different detection strategies have been proposed to

H. Li, B. Li, and J. Huang are with Guangdong Key Laboratory of Intelligent Information Processing and Shenzhen Key Laboratory of Media Security, Shenzhen University, Shenzhen 518060, China (e-mail: lihaodong, libin, jwhuang@szu.edu.cn).

Shunquan Tan is with College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China (tansq@szu.edu.cn).

B. Li is the correspondence author.

address these cases. Extensive experiments have been conducted to evaluate the performance of the proposed method. The experimental results show that the proposed method can effectively distinguish the DNG images from real ones in most cases. The main contributions of this paper are summarized as follows.

- We have analyzed the disparities between DNG images and real images. By measuring the similarities of DNG images and real images in different 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. Besides, we have also observed obvious disparities by assembling the R, G, B components together. These analyses and observations motivate the proposal of color statistical features.
- We have proposed an effective feature set for DNG image detection. The feature set consists of co-occurrence matrices
  extracted from the image high-pass filtering residuals of several color components. In order to make the feature dimension
  compact, binarization or truncation to the residuals is applied, and the elements of co-occurrence matrices are combined
  based on symmetric property. The proposed feature set is of low dimension, and achieves good detection performance
  even under the case of a small training set.
- We have designed different detection strategies based on whether the DNG images or the generative models are available. According to the level of available information, we have divided the detection situations into three cases, i.e., sample-aware, model-aware, and model-unaware. We have correspondingly designed the detection strategies, which employ different training data and classifiers. Specifically, binary classification is applied when DNG images or generated models are known, and one-class classification is utilized when the generative models are unknown. The effectiveness of the detection strategies has been validated by experiments.

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

## II. RELATED WORKS

## A. Image Generative Models

Given some training data, generative models can be trained to generate 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 any plausible samples similar to those coming from real world. Currently the most popular generative models based on deep neural networks include Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), and autoregressive models. Different from the former two, autoregressive models, such as PixelRNN [20] and PixelCNN [21], directly model the conditional distribution of pixels. As autoregressive models generate an image pixel by pixel, they spend more time on the generation process compared to GANs and VAEs. Moreover, autoregressive models usually produce images of poor quality, which would limit their applications. Hence, we only consider the generative models based on GANs and VAEs in this paper.

- 1) Generative Adversarial Networks: GAN was first proposed by Goodfellow et al. [22]. 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 a 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 eventually reproduces the distribution of real data, and the discriminator fails to distinguish between generated samples and real samples. Recently, many works [17]–[19], [23]–[25] have been proposed to improve the vanilla GAN. For example, Radford et al. [23] designed deep convolutional GAN (DCGAN), Arjovsky et al. [17] adopted Wasserstein distance in GAN to make the training more stable, Gulrajani et al. [18] made an improvement for Wasserstein GAN with gradient penalty (WGAN-GP), Karras et al. [19] proposed progressive growing of GAN (PGGAN) to improve the quality and variation of generated images.
- 2) Variational Autoencoders: Just like conventional autoencoders, VAE [26] has an encoder network and a decoder network. The encoder maps an image to a latent vector, while the decoder translate the latent vector back to an image. The major difference between VAEs and conventional autoencoders is that VAEs put constrains on the latent vector. Two kinds of loss are defined during training for optimization: one is the reconstruction loss of the image, and the other is the KL divergence loss of the latent variables. The variables of the latent vector are usually assumed to follow independent multivariate Gaussian distribution. After training, the encoder can represent a real sample as a vector of Gaussian variables, and the decoder can construct an realistic image when feeding it with a Gaussian vector, acting as a generator. The generative problem of VAEs can be formalized in the framework of probabilistic graphical models and it is easy to train. Different from GANs, VAEs do not have a discriminator that forces the generator to learn real data distribution. Since VAEs tend to generate blur images, some improved methods have been developed. Kingma et al. [27] combined VAE and GAN, and used the learned features of GAN to measure the reconstruction loss of VAE for obtaining better visual fidelity. Larsen et al. [28] tried to improve the visual

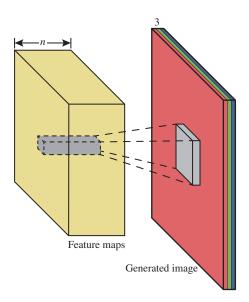


Fig. 1. The typical structure of the last layer of an image generator.

quality with inverse autoregressive flow. Hou et al. [29] added a deep feature consistent loss in VAE (DFC-VAE) to make the model produce better results.

In this paper, we consider the DNG images generated by DFC-VAE [29], DCGAN [23], WGAN-GP [18], and PGGAN [19].

## B. Face Spoofing Detection

Face spoofing attack attempts to bypass a face biometric system by presenting a fake face to the system [3]. It can be categorized into printing attack, replay attack, and 3D mask attack. In printing and replay attacks, images are printed/displayed and then recaptured by the authentication system, as gone through additional processing pipeline. Therefore, from the perspective of a detector, fake images are different from real images in visual quality because of the distortion introduced by recapturing. Thus, the fake faces can be detected via image quality assessment [30] or distortion analysis [31]. Since the distortions result in the changes of frequency spectrum and texture property, anti-spoofing detection can also be performed based on analysis of Fourier spectra [32] or textures [9]. DNG images do not go through the recapture pipeline, and thus have no recapture distortion. However, it is interesting to investigate whether the features designed for face spoofing detection are useful for DNG image detection. We will present the results in Section IV.

## III. DNG IMAGE DETECTION

In this section, we first analyze some possible artifacts of the DNG images and find out that some disparities exist between DNG images and real images in different color spaces. Then, we construct a feature set to capture the artifacts of DNG images so as to detect them. Finally, we discuss several detection scenarios and the corresponding detection strategies.

# A. Investigating DNG Images from the Perspective of Color

1) The generation pipeline of DNG images: In order to differentiate DNG images from real images, we investigate whether there are artifacts left during the generation of DNG images. Typically, an image generator takes a random latent vector as input, and employs several transpose convolutional<sup>1</sup> layers to gradually expand the spatial size of the random vector to produce an image. As shown in Fig. 1, in the last layer of the generator, several feature maps are transformed into a tensor with three channels via convolution, where the three channels respectively represent the red, green, and blue components of the generated image. As the image is generated in RGB space, the generator tends to learn the properties of real images in RGB space, while paying less attention to the properties in other color spaces. In this way, although the DNG images may look like real ones in RGB color space, there may be some differences in other color spaces. Moreover, a real image is captured from real scene, meaning that the color components are decomposed and digitalized from real world, while the color components in a DNG image are computed by three groups of convolutional weights without putting explicit constrains on their relations. Therefore, it is reasonable to assume that some inherent relations among the color components of the DNG images are different from real ones. In the following analysis, we will show some experimental evidences to support such an assumption.

<sup>&</sup>lt;sup>1</sup>Transpose convolution is also named deconvolution or fractional strided convolution in literatures.

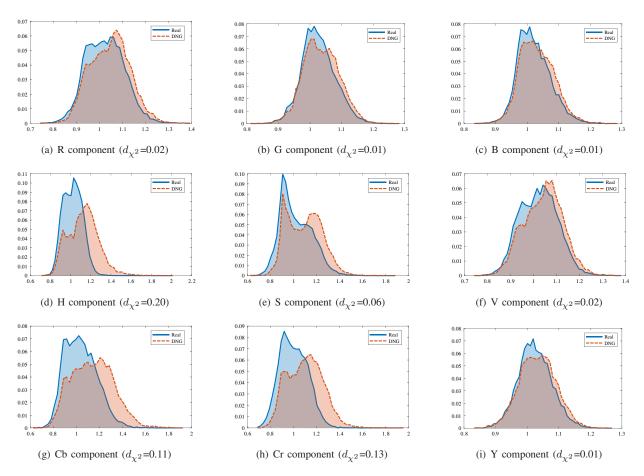


Fig. 2. The histograms  $H_{\mathrm{DNG}}^c$  (red) and  $H_{\mathrm{Real}}^c$  (blue) for different color components. The  $d_{\chi^2}(H_{\mathrm{DNG}}^c, H_{\mathrm{Real}}^c)$  values are included in the sub-captions.

2) Discernibility of different color components: In image processing, the RGB color space is widely used in cameras, software, and monitors for acquiring, representing, and displaying color. However, since the three color components are highly correlated with each other, as well as the fact that luminance and chrominance information are not well separated in the RGB space, it is beneficial to process the images in other color spaces, such as HSV and YCbCr.

In the following, we analyze the discernibility of three different color spaces, *i.e.*, RGB, HSV, and YCbCr, in distinguishing between DNG images and real images through analytical experiments. We try to use a metric to examine which color component is more discernible. To this aim, we first obtain the image statistics from different color components. Then we use a metric to evaluate the distance between the statistics from DNG images and those from real images. The larger the distance, the more discernible the color component.

We construct the image statistics as follows.

- a) For each color component, extract normalized histograms from some DNG images and real images. Compute the mean histograms by averaging the histograms from these two classes. Denote the mean histograms as  $\widetilde{H}^c$  and  $\overline{H}^c$  for DNG images and real images, respectively, where  $c \in \{R, G, B, H, S, V, Y, Cb, Cr\}$  represents different color components.
- b) Denote the histogram of the i-th image as  $H_i^c$ . We define a quantity called similarity index (SI) as

$$\lambda_i^c = \frac{d_{\chi^2}(H_i^c, \overline{H}^c)}{d_{\chi^2}(H_i^c, \widetilde{H}^c)},\tag{1}$$

where  $d_{\chi^2}(H_p, H_q)$  is the Chi-square distance for evaluating the similarity between two histograms  $H_p(x)$  and  $H_q(x)$  (x is the bin index) as

$$d_{\chi^2}(H_p, H_q) = \frac{1}{2} \sum_x \frac{(H_p(x) - H_q(x))^2}{H_p(x) + H_q(x)}.$$
 (2)

c) Compute the histogram of  $\lambda_i^c$  as the image statistics.

Denote the histograms of  $\{\lambda_i^c|i\in \mathrm{DNG}\}$  and  $\{\lambda_i^c|i\in \mathrm{Real}\}$  as  $H^c_{\mathrm{DNG}}$  and  $H^c_{\mathrm{Real}}$ . Compute the Chi-square distance  $d_{\chi^2}(H^c_{\mathrm{DNG}},H^c_{\mathrm{Real}})$  according to (2) as the discernible metric. It is expected that the large the distance, the better the discernibility.

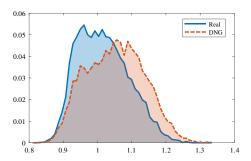


Fig. 3. The histograms  $H_{\rm DNG}^{RGB}$  (red) and  $H_{\rm Real}^{RGB}$  (blue).  $(d_{\chi^2}$ =0.05).

In our analytical experiments, we first use the CelebFaces Attributes dataset (CelebA) [33] to train an image generative model based on the WGAN-GP [18] method. With the trained model, we generate a large amount of DNG images. We randomly select 10000 DNG images and 10000 real images from the generated dataset and the CelebA dataset, respectively, and then compute the mean histograms,  $\widetilde{H}^c$  and  $\overline{H}^c$  ( $c \in \{R, G, B, H, S, V, Y, Cb, Cr\}$ ), as introduced above. Subsequently, we randomly select another 10000 DNG images and 10000 real images, and compute the SIs (i.e.,  $\lambda_i^c$ ) and obtain the histograms (i.e.,  $H_{\text{DNG}}^c$  and  $H_{\text{Real}}^c$ ). The histograms for different color components are shown in Fig. 2. From this figure, it can be observed that the overlapping regions of  $H_{\text{DNG}}^c$  and  $H_{\text{Real}}^c$  in the H, S, Cb, and Cr components are smaller than those in R, G, B, V, and Y components, implying that the disparities between DNG images and real images tend to be more obvious in the chrominance components.

The resultant discernible metrics, i.e.,  $d_{\chi^2}(H^c_{\rm DNG}, H^c_{\rm Real})$ , are also shown in the sub-captions of the corresponding sub-figures in Fig. 2. We can observe that the values of  $d_{\chi^2}(H^c_{\rm DNG}, H^c_{\rm Real})$  for the four chrominance components (i.e., H, S, Cb, and Cr), all being greater than 0.06, are larger than those for the luminance components (i.e., V and Y) and R, G, and B components, which are less than 0.02. These results indicate that the chrominance components are more discernible than the other components, implying that some statistical feature extracted from these chrominance components would be more effective in distinguishing between DNG images and real images.

3) Increasing discernibility by assembling RGB components: In previous part, we treated the R, G, and B components in RGB space separately. However, please note that the R, G, and B components of an image may be correlated with each other to some extents, which is quite different from the components in HSV and YCbCr spaces. Hence, analyzing the R, G, and B components individually may lead to the loss of discernibility. In order to increase the discernibility of the color information in RGB space, we try to assemble the R, G, and B components as a whole.

To this end, we regard R, G, and B values of each pixel within an image as a three-element tuple. Considering that a 24-bit RGB image (8-bit for each channel) will result in a histogram with a huge number of bins (i.e.,  $2^{24}$ ), we reduce the dimension of the histogram by uniformly quantizing each channel into 3-bit. In this way, a histogram with  $2^9 = 512$  bins can be obtained. As done in the previous part, we compute the histograms  $H_{\rm RGB}^{RGB}$  and  $H_{\rm Real}^{RGB}$ , and show them in Fig. 3. We can observe that the overlapping region of the two histograms are obviously smaller than those in the individual R, G, and B components as shown in Fig. 2(a)-2(c).

By calculating the Chi-square distances of the two histograms, we obtain  $d_{\chi^2}(H_{\rm DNG}^{RGB},H_{\rm Real}^{RGB})$ =0.05, which are larger than those for the individual R, G, and B components, meaning that there are sufficient disparities between DNG images and real images when assembling the R, G, and B components. Therefore, in addition to extracting features from the chrominance components, it is also helpful for detection if the features are extracted by considering the R, G, and B components together.

We have also conducted similar experiments by assembling H, S, V (or Y, Cb, Cr) components. The observed differences between DNG images and real images after assembling are less significant than those in individual chrominance components. It may be due to that the components in HSV and YCbCr are relatively de-correlated.

## B. Exacting Features from Color Components

From the previous subsection, we can conclude that there are some disparities between DNG images and real images from the perspective of color. Therefore, we would like to extract discernible features from color components. The overall framework of the proposed method is illustrated in Fig. 4. 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 this subsection, we describe the details of feature extraction process, and later we will discuss the designs and applications of classifier in the next subsection. In the feature extraction stage, we first perform high-pass filtering on images to suppress image contents. Then we compute the co-occurrence matrix on the high-pass filtering residuals, which are pre-processed by quantization or truncation for reducing the dimension of the co-occurrence matrix. The pre-processing

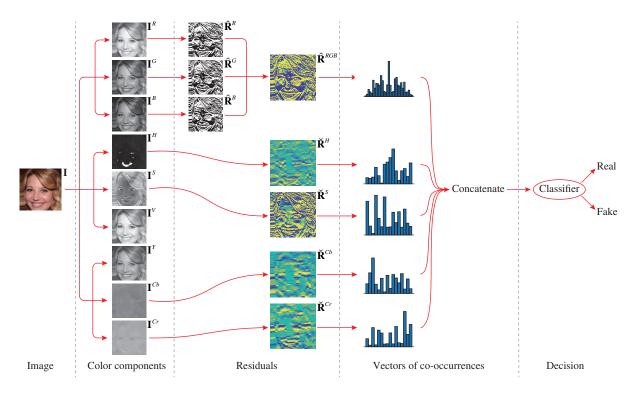


Fig. 4. The overall framework of the proposed method.

operation depends on which color component the residuals are coming from. Finally, the extracted co-occurrence matrix are merged to form a feature set.

1) Suppressing image contents: We plan to employ the co-occurrence matrix [34], which was widely used in image textural analysis as a kind of feature descriptor, to capture the disparities between DNG images and real images in different color components. However, instead of using the conventional gray level co-occurrence matrix of the image pixels, we propose to compute the co-occurrence matrix on image high-pass filtering residuals. The reason is that the contents of DNG images and real images are quite similar visually, especially in low-frequency representation such as the contour. As we know, human are less sensitive to high-frequency details. As a result, it is reasonable to suppress image contents so as to enhance high-frequency disparities. By using high-pass filtering, the high-frequency details can be well captured. In fact, extracting features from image high-pass filtering residuals have been successfully used in some applications, such as image steganalysis [35] and image forensics [36], to discover weak traces.

Given an image, denote its representations in RGB, HSV, and YCbCr spaces as  $I^{RGB}$ ,  $I^{HSV}$ , and  $I^{YCbCr}$ , respectively. For a color component  $I^c$ , its residual  $R^c$  can be obtained by:

$$\mathbf{R}^c = f(\mathbf{I}^c), c \in \{R, G, B, H, S, Cb, Cr\},$$
 (3)

where  $f(\cdot)$  is a high-pass filtering operation performed in spatial domain.

2) Binarization of RGB residuals: Based on the analysis in Section III-A, the R, G, B components should be treated as a whole for better discernibility. However, if there are many distinct element values, it would result in a co-occurrence matrix with a huge number of bins. As a result, we first binarize the residual image of each color component  $\mathbf{R}^c$  ( $c \in \{R, G, B\}$ ) by

$$\hat{\mathbf{R}}^c(x,y) = \begin{cases} 1, & \mathbf{R}^c(x,y) > 0, \\ 0, & \mathbf{R}^c(x,y) \le 0, \end{cases}$$
(4)

where (x,y) is the position index of an residual image element. Then, we obtain an assembled residual image  $\hat{\mathbf{R}}^{RGB}$  by

$$\hat{\mathbf{R}}^{RGB} = \hat{\mathbf{R}}^R \cdot 2^0 + \hat{\mathbf{R}}^G \cdot 2^1 + \hat{\mathbf{R}}^B \cdot 2^2. \tag{5}$$

In this way, the elements in  $\hat{\mathbf{R}}^{RGB}$  are within the range of [0,7].  $\hat{\mathbf{R}}^{RGB}$  is later used to compute the co-occurrence matrix.

$$\begin{pmatrix} 1 \\ -1 \end{pmatrix}, \quad \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

Fig. 5. The high-pass filters used for computing image residuals.

3) Truncation of chrominance residuals: Since H, S, Cb, Cr components represent different chrominance information of an image and thus have few correlations, we process them independently. In order to reduce the number of distinct values, the image residuals are truncated as follow:

$$\check{\mathbf{R}}^{c}(x,y) = \begin{cases} \tau, & \mathbf{R}^{c}(x,y) \ge \tau, \\ \mathbf{R}^{c}(x,y), & -\tau < \mathbf{R}^{c}(x,y) < \tau, \\ -\tau, & \mathbf{R}^{c}(x,y) \le -\tau, \end{cases}$$
(6)

where  $c \in \{H, S, Cb, Cr\}$  and  $\tau$  is the truncation threshold. The truncated image residuals are then used to compute the co-occurrence matrices.

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

$$\mathbf{C}(v_1, v_2, \dots, v_d) = \frac{1}{N} \sum_{x,y} \mathbb{1} \Big( \mathbf{V}(x, y) = v_1,$$

$$\mathbf{V}(x + \Delta x, y + \Delta y) = v_2, \dots,$$

$$\mathbf{V}(x + (d - 1)\Delta x, y + (d - 1)\Delta y) = v_d \Big),$$

$$(7)$$

where  $\mathbb{1}(\cdot)$  is an indicator function,  $(v_1, v_2, \dots, v_d)$  is the index of co-occurrence matrix, d is the order of co-occurrence matrix, N is the normalization factor, and  $\Delta x$ ,  $\Delta y$  are the offsets for two neighboring elements. The dimension of the co-occurrence matrix for  $\hat{\mathbf{R}}^{RGB}$  is  $8^d$ , while the dimensions for  $\check{\mathbf{R}}^H$ ,  $\check{\mathbf{R}}^S$ ,  $\check{\mathbf{R}}^{Cb}$ , and  $\check{\mathbf{R}}^{Cr}$  are  $(2\tau+1)^d$ .

Since the co-occurrence matrix is usually symmetric, we can decrease the feature dimension by combining the two bins,  $\mathbf{C}(v_1,v_2,\ldots,v_d)$  and  $\mathbf{C}(v_d,v_{d-1},\ldots,v_1)$ , into one bin. After combination, the dimensionality of co-occurrence matrix is substantially decreased: the co-occurrence matrix for  $\hat{\mathbf{R}}^{RGB}$  has only  $(8^d+8^{d-1})/2$  bins, and the co-occurrence matrices for  $\check{\mathbf{R}}^H$ ,  $\check{\mathbf{R}}^S$ ,  $\check{\mathbf{R}}^{Cb}$ , and  $\check{\mathbf{R}}^{Cr}$  have  $((2\tau+1)^d+(2\tau+1)^{d-1})/2$  bins. We 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 and improve the robustness to image translation.

5) Practical implementation: For the proposed scheme, one can balance the feature dimension and model effectiveness by using different high-pass filters and different parameters. In our practical implementation of extracting the proposed features, we simply use two difference operators, which are shown in Fig. 5, as high-pass filters to obtain the image residuals. The residuals are then processed as described above, where the truncation thresholds for  $\mathbf{R}^c(c \in \{H, S, Cb, Cr\})$  are set to  $\tau = 2$ . We set the order of co-occurrence matrix as d = 3 and the offsets as  $(\Delta x, \Delta y) \in \{(0,1), (1,0)\}$ . Therefore, we have 4 co-occurrence matrices (2 residuals  $\times$  2 offsets) for each color component, and we finally take the element-wise mean of the 4 co-occurrence matrices as the features. In total, a 588-D feature is extracted from each image, where the feature dimension is  $(8^3 + 8^2)/2 = 288$  for  $\hat{\mathbf{R}}^{RGB}$ , while the feature dimension is  $(5^3 + 5^2)/2 = 75$  for each of  $\hat{\mathbf{R}}^H$ ,  $\hat{\mathbf{R}}^S$ ,  $\hat{\mathbf{R}}^{Cb}$ , and  $\hat{\mathbf{R}}^{Cr}$ .

#### C. Detection Strategies

In practical applications, there are many kinds of generative models, and such models may be trained with different real image sources. As a result, DNG images generated by different models which are trained with different datasets may more or less exhibit different characteristics, leading to difficulties in distinguishing them from real images. Based on the information that an investigator can access, we divide the detection scenarios into three cases: sample-aware, model-aware, and model-unaware. These scenarios and the corresponding detection strategies are discussed as follows.

- 1) Sample-aware detection: In this case, the investigator can obtain some DNG images from a known generative model. This is the most simple case for the investigator. To perform the detection, the investigator can train a binary classifier with real images and DNG images, and uses the trained classifier to predict the class labels for the given images.
- 2) Model-aware detection: In this case, the investigator may know the generative process of DNG images, but he/she does not have any training images generated by the corresponding model, and has no idea about the real image dataset that was used to train the model. To perform the detection, the investigator needs to first train a generative model with the same network architecture by using a alternative dataset, and then use the trained model to produce DNG samples. With these samples, the investigator can train a binary classifier to detect DNG images. It is similar to the case of cross dataset validation in many applications. It is expected that better generalization ability brings better performance.

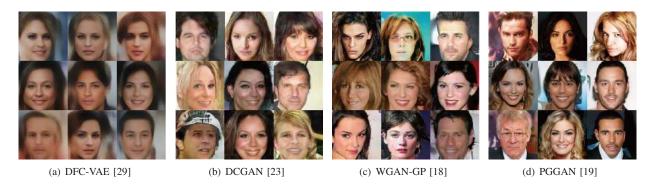


Fig. 6. Sample images generated by different methods.

3) Model-unaware detection: It is not rare that the investigator does not have any DNG image samples nor have any knowledge about the generative model. This is the most challenging case presented to the investigator. To cope with this case, the investigator may train a one-class classifier with real images and use the classifier to detect whether the testing image is real or DNG.

## IV. EXPERIMENTS

In this section, we will examine the performance of the proposed method. First we introduce the common experimental settings. Then we present the experimental results for the three kinds of detection scenarios as described in Section III-C.

## A. Experimental Setups

- 1) Real image datasets: In the experiments, we adopt three face image datasets for evaluating the performance of the proposed method. All images in these datasets are regarded as real images. The real image sets are denoted as  $\mathcal{R}$ , and we use subscripts to indicate the image source. The details of the datasets are as follows.
  - CelebFaces Attributes (CelebA) [33]: This dataset consists of more than 200K celebrity face images. We use the "Align & Cropped" PNG images in this dataset, and then crop the facial region with  $138 \times 138$  from each image to remove the background, and then resize the cropped region into  $64 \times 64$  and  $128 \times 128$ , respectively. The resulting datasets are denoted as  $\mathcal{R}_{C.64}$  and  $\mathcal{R}_{C.128}$ .
  - High-quality CelebA (HQ-CelebA) [19]: This dataset contains 30000 face images in 1024×1024 resolution, which are
    obtained by applying several processing to the in-the-wild images in the CelebA dataset. This dataset is denoted as R<sub>H-1024</sub>.
  - Labeled Faces in the Wild (LFW) [37], [38]: This dataset contains 13233 face images. We used the calibrated version [38] of this dataset, which contain images of size 250×250. We also crop the facial region with 150×150 from each image to remove the background, and then resize the cropped region into 128×128. The resulting dataset is denoted as  $\mathcal{R}_{L-128}$ .
- 2) Generative models and DNG image datasets: Four typical types of generative models, including DFC-VAE [29], DCGAN [23], WGAN-GP [18] and PGGAN [19], are used in our experiments. Among them, DFC-VAE is based on variational autoencoder while the others are all based on generative adversarial networks. For the first three models<sup>2</sup>, we adapt the network architectures to generate images with two sizes, *i.e.*,  $64 \times 64$  and  $128 \times 128$ . For PGGAN, we download 30000  $1024 \times 1024$  generated images shared online. For simplicity, in the following context we denote a set of DNG images as  $\mathcal{G}$ , and use a superscript to represent the type of generative model and a subscript to represent the real image dataset used in training the generative model. For example,  $\mathcal{G}_{C-64}^{DCGAN}$  denotes the image set generated by a DCGAN model trained with  $64 \times 64$  CelebA images. Fig. 6 shows some generated samples of different methods. It is observed that the images generated by DFC-VAE look more blurred than others, and some images generated by DCGAN may present unnatural textures. WGAN-GP significantly improves the visual quality of generated images compared to DFC-VAE and DCGAN, while PGGAN generates more visually pleasing images since it produces images with high resolution. We have invited 10 people to differentiate the images in  $\mathcal{R}_{H-1024}$  and  $\mathcal{G}_{H-1024}^{PCGAN}$ . The obtained average accuracy was about 80%, meaning that the images generated by PGGAN are not easy to be identified.

#### B. Sample-aware Detection

In this subsection, we evaluate the performance of sample-aware detection by using binary classifier trained with known DNG image samples. Seven DNG image sets are used in the experiments, *i.e.*,  $\mathcal{G}_{\text{C-64}}^{\text{DFC-VAE}}$ ,  $\mathcal{G}_{\text{C-64}}^{\text{DCGAN}}$ ,  $\mathcal{G}_{\text{C-128}}^{\text{DCGAN}}$ ,  $\mathcal{G}_{\text{C-128}}^{\text{DCGAN}}$ ,  $\mathcal{G}_{\text{C-128}}^{\text{DCGAN}}$ ,  $\mathcal{G}_{\text{C-128}}^{\text{DCGAN}}$ , and  $\mathcal{G}_{\text{H-1024}}^{\text{PGGAN}}$ . The number of images in each of these datasets is the same as that in the corresponding real image datasets. We

<sup>&</sup>lt;sup>2</sup>We use the implementation of DFC-VAE at: https://www.github.com/houxianxu/DFC-VAE and the implementation of WGAN-GP and DCGAN at: https://www.github.com/igul222/improved\_wgan\_training.

 $\label{table I} \textbf{TABLE I}$  Detection results for full-frame face images under sample-aware scenario.

| Detector  | Training time (s) | Testing set   | FPR (%) | FNR (%) | ACC (%) |
|---|-------------------|---|---------|---------|---------|
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{DFC-VAE}} ight)$        | 64.0              | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{DFC-VAE}}\}$      | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}} ight)$          | 91.2              | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}}\}$        | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}} ight)$        | 91.3              | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}}\}$      | 0.17    | 0.05    | 99.89   |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DFC-VAE}} ight)$      | 49.0              | $\{\mathcal{R}_{	ext{C-}128}, \mathcal{G}_{	ext{C-}128}^{	ext{DFC-VAE}}\}$    | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}} ight)$        | 46.3              | $\{\mathcal{R}_{	ext{C-}128}, \mathcal{G}_{	ext{C-}128}^{	ext{DCGAN}}\}$      | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}} ight)$      | 77.0              | $\{\mathcal{R}_{	ext{C-}128}, \mathcal{G}_{	ext{C-}128}^{	ext{WGAN-GP}}\}$    | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{H-1024}},\mathcal{G}_{	ext{H-1024}}^{	ext{PGGAN}} ight)$      | 13.8              | $\{\mathcal{R}_{\text{H-1024}},\mathcal{G}_{\text{H-1024}}^{PGGAN}\}$         | 0.01    | 0.00    | 100.00  |
| $\phi_{\text{CoALBP+LPQ}}\left(\mathcal{R}_{\text{C-64}},\mathcal{G}_{\text{C-64}}^{\text{DFC-VAE}}\right)$ | 86.8              | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{DFC-VAE}}\}$      | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{CoAlbP+LPQ}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}} ight)$        | 1420.4            | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}}\}$        | 0.01    | 0.01    | 99.99   |
| $\phi_{	ext{CoAlbP+LPQ}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}} ight)$      | 7882.8            | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}}\}$      | 0.12    | 0.01    | 99.93   |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DFC-VAE}} ight)$    | 66.3              | $\{\mathcal{R}_{	ext{C-}128}, \mathcal{G}_{	ext{C-}128}^{	ext{DFC-VAE}}\}$    | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}} ight)$      | 73.2              | $\{\mathcal{R}_{	ext{C-}128}, \mathcal{G}_{	ext{C-}128}^{	ext{DCGAN}}\}$      | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{Coalbp+lpq}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}} ight)$    | 1555.9            | $\{\mathcal{R}_{\text{C-128}}, \mathcal{G}_{\text{C-128}}^{\text{WGAN-GP}}\}$ | 0.01    | 0.01    | 99.99   |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{H-1024}},\mathcal{G}_{	ext{H-1024}}^{	ext{PGGAN}} ight)$    | 13.2              | $\{\mathcal{R}_{\text{H-1024}},\mathcal{G}^{\text{PGGAN}}_{\text{H-1024}}\}$  | 0.01    | 0.00    | 100.00  |
| $\phi_{	ext{DCGAN}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}} ight)$             | -                 | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}}\}$        | 0.33    | 0.52    | 99.57   |
| $\phi_{	ext{DCGAN}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}} ight)$           | -                 | $\{\mathcal{R}_{	ext{C-}128}, \mathcal{G}_{	ext{C-}128}^{	ext{DCGAN}}\}$      | 1.25    | 1.04    | 98.86   |
| $\phi_{	ext{WGAN-GP}}\left(\mathcal{R}_{	ext{C-64}},\mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}} ight)$         | -                 | $\{\mathcal{R}_{	ext{C-64}}, \mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}}\}$      | 42.28   | 37.51   | 60.10   |
| $\phi_{	ext{WGAN-GP}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}} ight)$       | -                 | $\{\mathcal{R}_{\text{C-128}},\mathcal{G}_{\text{C-128}}^{\text{WGAN-GP}}\}$  | 35.78   | 38.67   | 62.77   |

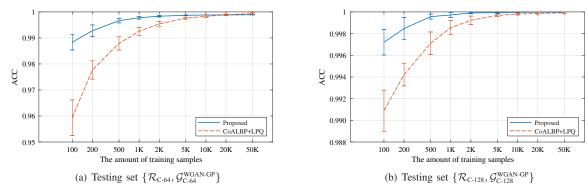


Fig. 7. Detection results with different amounts of training samples under sample-aware scenario.

train 7 binary classifiers by regarding the DNG images as positive samples and the corresponding real images as negative samples. An ensemble of LDA (linear discriminative analysis) base learners [39] is used as the classifier. For comparison, we train another 7 classifiers with a feature set proposed in [9], which is designed for face spoofing detection. The feature set is composed of Co-Occurrence of Adjacent Local Binary Patterns (CoALBP) features and Local Phase Quantization (LPQ) features. The dimension of CoALBP+LPQ is 19968-D. In the default setting, 25% of the real and the DNG images from the respective dataset are randomly selected and used in the training stage, while the remaining 75% images are used for testing. We repeat the training and testing stages 10 times by independently splitting the training and testing samples and compute the mean testing results. The performance is measured by false positive rate (FPR), false negative rate (FNR), and overall accuracy (ACC).

The performance of our method and that of the CoALBP+LPQ method are reported in Table I, where  $\phi_{\text{Proposed}}(\mathcal{R},\mathcal{G})$  and  $\phi_{\text{CoALBP+LPQ}}(\mathcal{R},\mathcal{G})$  respectively denote the classifier trained with the proposed features and that with the CoALBP+LPQ features, using the training real image set  $\mathcal{R}$  and DNG image set  $\mathcal{G}$ . From Table I we can observe that the detection accuracies for all the cases are higher than 99%, meaning that both the proposed feature and the CoALBP+LPQ feature can achieve quite good detection performance. Only some occasional errors are made when performing classification on  $\{\mathcal{R}_{\text{C-64}},\mathcal{G}^{\text{WGAN-GP}}_{\text{C-64}}\}$ . The good results indicate that the DNG images can be easily identified when we have a targeted detector that is trained with real images and the corresponding DNG images. Comparing with the high-dimensional CoALBP+LPQ feature, the proposed feature set is of low dimension and it can achieve almost the same detection performance. Besides, the proposed method significantly reduces the computation cost due to the low feature dimension. In our implementation, the time for extracting the proposed feature is about 1/10 of that for extracting the CoALBP+LPQ feature. As shown in Table I, the time for training classifiers

TABLE II
DETECTION RESULTS FOR THE MODEL-AWARE SCENARIO.

| Detector   | Testing set  | FPR (%) | FNR (%) | ACC (%) |
|--|--|---------|---------|---------|
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}} ight)$     |  | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}} ight)$   | $\{\mathcal{R}_{	ext{L-128}}, \mathcal{G}_{	ext{L-128}}^{	ext{DCGAN}}\}$   | 0.00    | 0.00    | 100.00  |
| $\phi_{	ext{DCGAN}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}} ight)$        |  | 21.49   | 21.28   | 78.62   |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}} ight)$   |  | 0.00    | 0.16    | 99.92   |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}} ight)$ | $\{\mathcal{R}_{	ext{L-}128}, \mathcal{G}_{	ext{L-}128}^{	ext{WGAN-GP}}\}$ | 0.00    | 0.07    | 99.97   |
| $\phi_{	ext{WGAN-GP}}\left(\mathcal{R}_{	ext{C-128}},\mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}} ight)$    |  | 41.72   | 38.57   | 59.86   |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{L-128}},\mathcal{G}_{	ext{L-128}}^{	ext{DCGAN}} ight)$     |  | 0.54    | 0.00    | 99.73   |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{L-128}},\mathcal{G}_{	ext{L-128}}^{	ext{DCGAN}} ight)$   | $\{\mathcal{R}_{	ext{C-128}}, \mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}}\}$   | 0.63    | 0.00    | 99.69   |
| $\phi_{	ext{DCGAN}}\left(\mathcal{R}_{	ext{L-128}},\mathcal{G}_{	ext{L-128}}^{	ext{DCGAN}} ight)$        |  | 53.84   | 16.64   | 64.76   |
| $\phi_{	ext{Proposed}}\left(\mathcal{R}_{	ext{L-128}},\mathcal{G}_{	ext{L-128}}^{	ext{WGAN-GP}} ight)$   |  | 3.39    | 0.05    | 98.28   |
| $\phi_{	ext{CoALBP+LPQ}}\left(\mathcal{R}_{	ext{L-}128},\mathcal{G}_{	ext{L-}128}^{	ext{WGAN-GP}} ight)$ | $\{\mathcal{R}_{	ext{C-128}}, \mathcal{G}_{	ext{C-128}}^{	ext{WGAN-GP}}\}$ | 15.64   | 0.00    | 92.18   |
| $\phi_{	ext{WGAN-GP}}\left(\mathcal{R}_{	ext{L-128}},\mathcal{G}_{	ext{L-128}}^{	ext{WGAN-GP}} ight)$    |  | 70.04   | 27.09   | 51.44   |

with the proposed feature is usually shorter than that for training with CoALBP+LPQ. Note that  $\phi_{\text{CoALBP+LPQ}}$  ( $\mathcal{R}_{\text{C-64}}$ ,  $\mathcal{G}_{\text{C-64}}^{\text{DCGAN}}$ ),  $\phi_{\text{CoALBP+LPQ}}$  ( $\mathcal{R}_{\text{C-64}}$ ,  $\mathcal{G}_{\text{C-64}}^{\text{WGAN-GP}}$ ), and  $\phi_{\text{CoALBP+LPQ}}$  ( $\mathcal{R}_{\text{C-128}}$ ,  $\mathcal{G}_{\text{C-128}}^{\text{WGAN-GP}}$ ) require more time due to the search of the best number of feature dimension in base learners.

Since the discriminators of GANs aim to differentiate between generated images and real images, we have also examined the detection performance of the discriminators of trained models of DCGAN and WGAN-GP, denoted by  $\phi_{\text{DCGAN}}(\mathcal{R}, \mathcal{G})$  and  $\phi_{\text{WGAN-GP}}(\mathcal{R}, \mathcal{G})$ , respectively. The detection results are shown in the last part of Table I, from which we observe that the discriminator of DCGAN performs well (ACC > 98%) while the discriminator of WGAN-GP performs poorly ( $ACC \approx 60\%$ ). The result is not surprising for that WGAN-GP leads to better visual quality and better undetectability than DCGAN.

In order to investigate how the performance is affected by the amount of training samples, we reduce the size of the training set. Fig. 7 shows the detection accuracies of the proposed method and those of CoALBP+LPQ on  $\{\mathcal{R}_{\text{C-64}}, \mathcal{G}_{\text{C-64}}^{\text{WGAN-GP}}\}$  and  $\{\mathcal{R}_{\text{C-128}}, \mathcal{G}_{\text{C-128}}^{\text{WGAN-GP}}\}$  with different amounts of training samples. The accuracies are obtained by repeated training and testing 10 times, and the standard deviations of accuracies are shown as error bars. It can be observed that the proposed method significantly outperforms CoALBP+LPQ when the amount of training samples is small. The testing performance of the proposed method tends to be stable when the training samples are more than 2000, while that of CoALBP+LPQ becomes stable when the training samples are 20000. It can also be observed that the proposed method has lower standard deviations, meaning that its performance is more robust.

## C. Model-aware Detection

It is reasonable to assume that the architecture of generative model is known, but the real images used to train the model are not available to the investigator. This is the scenarios we called model-aware detection. In this subsection, we evaluate the performance of the proposed method under this situation. To this end, we perform experiments as follows. First, by using two distinct real image sets, we generate two DNG image sets independently by using a generative model. Then we train a classifier upon one real image set and its corresponding DNG image set, and test its performance on another real image set and its corresponding DNG image set. We use DCGAN and WGAN-GP as the generative models in our experiments, and use  $\mathcal{R}_{\text{C-128}}$  and  $\mathcal{R}_{\text{L-128}}$  as the real image sets.

We compare the proposed method with CoALBP+LPQ and the discriminators of GANs. The detection results are shown in Table II. It can be observed that the proposed method and the CoALBP+LPQ method, which are based on hand-crafted features, work much better than the GAN discriminators. It implies that the GAN discriminators are not suitable for the model-aware detection, for that the discriminators are overfitted to the training data and have poor generalization performance. It can also be observed that the proposed method and the CoALBP+LPQ method have similar performance for the first three cases, while the proposed method outperforms the CoALBP+LPQ method for the case when  $\{\mathcal{R}_{\text{L-128}}, \mathcal{G}_{\text{L-128}}^{\text{WGAN-GP}}\}$  is used for training and  $\{\mathcal{R}_{\text{C-128}}, \mathcal{G}_{\text{L-128}}^{\text{WGAN-GP}}\}$  is used for testing.

## D. Model-unaware Detection

In a more practical situation, both the DNG images and the corresponding generated model are not available to the investigator. This is the so-called model-unaware detection scenario. In this subsection, we evaluate the performance of proposed method under this situation. In order to perform the detection, we use some real images to build a one-class classifier, which fits a model to describe the distribution of real images and regard the DNG images as outliers. Hence, by feeding testing images to

TABLE III
DETECTION RESULTS (%) FOR MODEL-UNAWARE SCENARIO.

|   | $\{\mathcal{R}_{\text{C-64}}\}$  | $\{\mathcal{G}_{	ext{C-64}}^{	ext{DFC-VAE}}\}$  | $\{\mathcal{G}_{	ext{C-64}}^{	ext{DCGAN}}\}$  | $\{\mathcal{G}_{	ext{C-64}}^{	ext{WGAN-GP}}\}$  |
|---|----------------------------------|---|---|---|
| $arphi_{	ext{Proposed}}^{0.10}\left(\mathcal{R}_{	ext{C-64}} ight) \ arphi_{	ext{Proposed}}^{0.05}\left(\mathcal{R}_{	ext{C-64}} ight)$   | 89.59                            | 99.97   | 100.00  | 51.50   |
|   | 94.31                            | 99.51   | 100.00  | 32.04   |
|   | $\{\mathcal{R}_{\text{C-128}}\}$ | $\{\mathcal{G}_{	ext{C-128}}^{	ext{DFC-VAE}}\}$ | $\{\mathcal{G}_{	ext{C-128}}^{	ext{DCGAN}}\}$ | $\{\mathcal{G}^{	ext{WGAN-GP}}_{	ext{C-128}}\}$ |
| $arphi_{	ext{Proposed}}^{0.10}\left(\mathcal{R}_{	ext{C-}128} ight) \ arphi_{	ext{Proposed}}^{0.05}\left(\mathcal{R}_{	ext{C-}128} ight)$ | 89.40                            | 99.69   | 100.00  | 99.29   |
|   | 94.29                            | 96.62   | 100.00  | 96.46   |
|   | $\{\mathcal{R}_{L\text{-}128}\}$ | -   | $\{\mathcal{G}_{	ext{L-}128}^{	ext{DCGAN}}\}$ | $\{\mathcal{G}_{	ext{L-}128}^{	ext{WGAN-GP}}\}$ |
| $arphi_{	ext{Proposed}}^{0.10}\left(\mathcal{R}_{	ext{C-}128} ight) \ arphi_{	ext{Proposed}}^{0.05}\left(\mathcal{R}_{	ext{C-}128} ight)$ | 86.90                            | -   | 100.00  | 97.70   |
|   | 94.18                            | -   | 100.00  | 92.43   |

the one-class classifier, we can identify the DNG images once they are not predicted as real ones. The extracted features as well as the one-class classifier play an important role in the detection. In our experiment, we use the popular LIBSVM [40] as the one-class classifier. The Gaussian kernel is selected and its parameter  $\gamma$  is determined via a grid search. The parameter nu, which controls the upper bound of training error (regarding the training real images as outliers), is set as 0.10 and 0.05, respectively. The trained one-class classifier is denoted as  $\varphi^{nu}_{\text{Proposed}}(\mathcal{R})$ .

First, we randomly select 10000 real images from  $\mathcal{R}_{\text{C-64}}$  as training images, and perform the testing on the rest of real images in  $\mathcal{R}_{\text{C-64}}$  and the DNG images created by different generative models. The results are shown in the first part of Table III. On the one hand, it can be observed that the detection performance for real images is related to the parameter nu. Specifically, when nu equals 0.10 and 0.05, the testing errors for real images are correspondingly close to 10% and 5%, respectively. It means that the trained classifier can adequately model the distribution of images in the real image set. On the other hand, the one-class classifiers can accurately detect the images generated by DFC-VAE and DCGAN, while they cannot satisfactorily reveal the images generated by WGAN-GP. It implies that WGAN-GP generates more realistic images in this situation. Then, we replace the image size by  $128 \times 128$  and perform the same experiments on  $\mathcal{R}_{\text{C-128}}$  as above. The results are shown in the second part of Table III. It can be observed that the one-class classifiers can achieve good performance, even for the images generated by WGAN-GP. Finally, we use the classifiers trained by  $\mathcal{R}_{\text{C-128}}$  to test the real images in  $\mathcal{R}_{\text{L-128}}$  and the related DNG images  $^3$ . The results are shown in the third part of Table III. It can be observed that the testing accuracies for real images in  $\mathcal{R}_{\text{L-128}}$  are slightly lower than those for real images in  $\mathcal{R}_{\text{C-128}}$ , due to the fact that the testing and training samples are from different datasets. On the other hand, it is observed that perfect performance is still achieved for detecting the DNG images in  $\mathcal{G}_{\text{C-128}}^{\text{NGAN-GP}}$ , while the performance for the DNG images in  $\mathcal{G}_{\text{C-128}}^{\text{NGAN-GP}}$  is slightly degraded compared to that for the DNG images  $\mathcal{G}_{\text{C-128}}^{\text{NGAN-GP}}$ .

It takes about 20 minutes for the SVM to converge during the training stage for one single parameter  $\gamma$ , performing on a computer with Intel Xeon E5-2630 v2 CPU and 128GB RAM. We have also tried to use the CoALBP+LPQ features to train one-class SVM. However, since the feature is of high dimension, it is difficult for the support vector machine to learn the decision boundary. It shows the importance of using a compact feature set. Based on these experimental results, it is promising to apply a one-class classifier with the proposed features for detecting DNG images in the model-unaware situation.

## V. CONCLUSION

In this paper, we have investigated the issue of detecting deep network generated images. We have analyzed the disparities between DNG images and real images and obtained some useful observations. Based on the observations, a feature set based on color statistical features is proposed. The feature set is compact and effective. According to the availability of information in practice, three detection scenarios are considered and the corresponding detection strategies are designed. To evaluate the performance of the proposed method, extensive experiments have been conducted. The experimental results show that the proposed features equipped with a binary classifier can effectively differentiate between DNG images and real images when DNG samples or generative models are available. Moreover, the proposed features together with a one-class classifier can also achieve good performance when the generative models are unknown.

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 acceptable 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

 $<sup>^{3}</sup>$ We have trained DFC-VAE with  $\mathcal{R}_{L-128}$ , but the obtained model outputs random patterns rather than faces. Therefore, DFC-VAE is not included in this experiment.

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 the proposed method to meet requirements in practical detection scenarios, and extend the proposed method to detect DNG images with diverse contents. Furthermore, we will try to detect the modern image processing techniques that utilizes deep networks based generative models as backend, for example, image inpainting with GANs.

#### ACKNOWLEDGEMENT

We would like to thank the authors of [9], [18], [19], [23], [29] for sharing their codes and/or datasets online.

#### REFERENCES

- [1] M. C. Stamm, M. Wu, and K. R. Liu, "Information forensics: An overview of the first decade," IEEE Access, vol. 1, pp. 167-200, 2013.
- [2] P. Korus, "Digital image integrity-a survey of protection and verification techniques," Digital Signal Process., vol. 71, pp. 1-26, 2017.
- [3] J. Galbally, S. Marcel, and J. Fierrez, "Biometric antispoofing methods: A survey in face recognition," IEEE Access, vol. 2, pp. 1530-1552, 2014.
- [4] W. Luo, J. Huang, and G. Qiu, "JPEG error analysis and its applications to digital image forensics," *IEEE Trans. Inf. Forensics Security*, vol. 5, no. 3, pp. 480–491, 2010.
- [5] T. Bianchi and A. Piva, "Image forgery localization via block-grained analysis of JPEG artifacts," *IEEE Trans. Inf. Forensics Security*, vol. 7, no. 3, pp. 1003–1017, 2012.
- [6] P. Korus and J. Huang, "Multi-scale fusion for improved localization of malicious tampering in digital images," *IEEE Trans. Image Process.*, vol. 25, no. 3, pp. 1312–1326, 2016.
- [7] H. Li, W. Luo, X. Qiu, and J. Huang, "Image forgery localization via integrating tampering possibility maps," *IEEE Trans. Inf. Forensics Security*, vol. 12, no. 5, pp. 1240–1252, 2017.
- [8] K. Patel, H. Han, and A. K. Jain, "Secure face unlock: Spoof detection on smartphones," *IEEE Trans. Inf. Forensics Security*, vol. 11, no. 10, pp. 2268–2283, 2016.
- [9] Z. Boulkenafet, J. Komulainen, and A. Hadid, "Face spoofing detection using colour texture analysis," *IEEE Trans. Inf. Forensics Security*, vol. 11, no. 8, pp. 1818–1830, 2016.
- [10] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," arXiv preprint arXiv:1609.03499, 2016.
- [11] 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.
- [12] J.-Y. Zhu, T. Park, P. Isola, and 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.
- [13] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Proc. Conf. Neural Information Processing Systems (NIPS)*, 2017, pp. 700–708.
- [14] S. Iizuka, E. Simo-Serra, and H. Ishikawa, "Globally and locally consistent image completion," ACM Trans. Graphics, vol. 36, no. 4, pp. 107:1–107:14, 2017
- [15] J. Snow. (2017, Nov) AI could set us back 100 years when it comes to how we consume news. [Online]. Available: https://www.technologyreview.com/s/609358
- [16] W. Knight. (2018, May) The us military is funding an effort to catch deepfakes and other AI trickery. [Online]. Available: https://www.technologyreview.com/s/611146
- [17] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," arXiv preprint arXiv:1701.07875, 2017.
- [18] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of Wasserstein GANs," in *Proc. Conf. Neural Information Processing Systems (NIPS)*, 2017, pp. 5769–5779.
- [19] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of GANs for improved quality, stability, and variation," in *Proc. Int. Conf. Learning Representations (ICLR)*, 2018.
- [20] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," in Proc. Int. Conf. Machine Learning (ICML), 2016, pp. 1747–1756.
- [21] A. van den Oord, N. Kalchbrenner, L. Espeholt, k. kavukcuoglu, O. Vinyals, and A. Graves, "Conditional image generation with PixelCNN decoders," in *Proc. Conf. Neural Information Processing Systems (NIPS)*, 2016, pp. 4790–4798.
- [22] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. Conf. Neural Information Processing Systems (NIPS)*, 2014, pp. 2672–2680.
- [23] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," in *Proc. Int. Conf. Learning Representations (ICLR)*, 2016.
- [24] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. P. Smolley, "Least squares generative adversarial networks," in *Proc. IEEE Int. Conf. Computer Vision (ICCV)*, 2017, pp. 2813–2821.
- [25] J. Zhao, M. Mathieu, and Y. LeCun, "Energy-based generative adversarial network," in Proc. Int. Conf. Learning Representations (ICLR), 2017.
- [26] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *Proc. Int. Conf. Learning Representations (ICLR)*, 2014.
- [27] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling, "Improved variational inference with inverse autoregressive flow," in *Proc. Conf. Neural Information Processing Systems (NIPS)*, 2016, pp. 4743–4751.
- [28] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, "Autoencoding beyond pixels using a learned similarity metric," in *Proc. Int. Conf. Machine Learning (ICML)*, 2016, pp. 1558–1566.
- [29] X. Hou, L. Shen, K. Sun, and G. Qiu, "Deep feature consistent variational autoencoder," in *Proc. IEEE Winter Conf. Applications of Computer Vision (WACV)*, 2017, pp. 1133–1141.
- [30] J. Galbally and S. Marcel, "Face anti-spoofing based on general image quality assessment," in *Proc. Int. Conf. Pattern Recognition (ICPR)*, 2014, pp. 1173–1178.
- [31] D. Wen, H. Han, and A. K. Jain, "Face spoof detection with image distortion analysis," *IEEE Trans. Inf. Forensics Security*, vol. 10, no. 4, pp. 746–761, 2015.
- [32] J. Li, Y. Wang, T. Tan, and A. K. Jain, "Live face detection based on the analysis of fourier spectra," in Proc. SPIE, vol. 5404, 2004, pp. 296-304.
- [33] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in Proc. IEEE Int. Conf. Computer Vision (ICCV), 2015, pp. 3730–3738.
- [34] R. M. Haralick, K. Shanmugam, and I. Dinstein, "Textural features for image classification," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-3, no. 6, pp. 610–621, Nov. 1973.
- [35] J. Fridrich and J. Kodovsky, "Rich models for steganalysis of digital images," IEEE Trans. Inf. Forensics Security, vol. 7, no. 3, pp. 868-882, 2012.

- [36] H. Li, W. Luo, X. Qiu, and J. Huang, "Identification of various image operations using residual-based features," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 28, no. 1, pp. 31–45, Jan. 2018.
- [37] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled Faces in the Wild: A database for studying face recognition in unconstrained environments," University of Massachusetts, Amherst, Tech. Rep. 07-49, Oct. 2007.
- [38] G. B. Huang, M. Mattar, H. Lee, and E. Learned-Miller, "Learning to align from scratch," in *Proc. Conf. Neural Information Processing Systems (NIPS)*, 2012, pp. 764–772.
- [39] J. Kodovsky, J. Fridrich, and V. Holub, "Ensemble classifiers for steganalysis of digital media," *IEEE Trans. Inf. Forensics Security*, vol. 7, no. 2, pp. 432–444, Apr. 2012.
- [40] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," ACM Trans. Intell. Syst. Technol., vol. 2, pp. 27:1–27:27, 2011, software available at http://www.csie.ntu.edu.tw/~cjlin/libsvm.