

Exposing GAN-synthesized Faces Using Landmark Locations

Xin Yang*, Yuezun Li*, Honggang Qi[†], Siwei Lyu*

* Computer Science Department, University at Albany, State University of New York, USA

[†] School of Computer and Control Engineering, University of the Chinese Academy of Sciences, China

ABSTRACT

Generative adversary networks (GANs) have recently led to highly realistic image synthesis results. In this work, we describe a new method to expose GAN-synthesized images using the locations of the facial landmark points. Our method is based on the observations that the facial parts configuration generated by GAN models are different from those of the real faces, due to the lack of global constraints. We perform experiments demonstrating this phenomenon, and show that an SVM classifier trained using the locations of facial landmark points is sufficient to achieve good classification performance for GAN-synthesized faces.

KEYWORDS

Image Forensics, GANs, Facial landmarks

ACM Reference Format:

Xin Yang*, Yuezun Li*, Honggang Qi[†], Siwei Lyu*. 2019. Exposing GAN-synthesized Faces Using Landmark Locations. In . ACM, New York, NY, USA, 6 pages.

1 INTRODUCTION

The fast advancement of artificial intelligence technologies and the increasing availability of large volume of online images and videos and high-throughput computing hardware have revolutionized the tools to generate visually realistic images and videos. These technologies are becoming more efficient and accessible to more users. The recent developments in deep neural networks [1, 3, 4, 21], and in particular, the generative adversary networks (GANs) [5], have spawned a new type of image synthesis methods that can produce images with high levels of realism. Figure 1 shows a few examples of GAN synthesized faces, with very impressive results obtained using recent GAN-based methods [7, 8].

The increasing sophistication of GAN-synthesized images also has the negative effect of fake visual media, and the most damaging examples of which are perhaps the fabricated or manipulated human faces since faces carry the most identifiable information of a person. The wide spread of fake media with GAN-synthesized faces raise significant ethical, legal and security concerns, and there is an urgent need for methods that can detect GAN-synthesized faces in images and videos.

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

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Unlike previous image/video manipulation methods, realistic images are generated completely from random noise through a deep neural network. Current detection methods are based on low level features such as color disparities [10, 13], or using the whole image as input to a neural network to extract holistic features [19].

In this work, we develop a new GAN-synthesized face detection method based on a more semantically meaningful features, namely the locations of facial landmark points. This is because the GAN-synthesized faces exhibit certain abnormality in the facial landmark locations. Specifically, The GAN-based face synthesis algorithm can generate face parts (*e.g.*, eyes, nose, skin, and mouth, etc) with a great level of realistic details, yet it does not have an explicit constraint over the locations of these parts in a face. To make an analogy, the GAN-based face synthesis method works like players in a game of Fukuwarai¹, it has all the face parts, but lacks in placing them in a natural and coherent way as in a real face.

We show that these abnormalities in the configuration of facial parts in GAN-synthesized faces can be revealed using the locations of the facial landmark points (*e.g.*, tips of the eyes, nose and the mouth) automatically detected on faces. To accommodate the variations in shape, orientation and scale of different faces, we further normalize all the facial landmarks to the same standard coordinate system. We then used the normalized locations of these facial landmarks as features for a simple SVM classifier. The landmark location based SVM classifier is tested on faces generated with the state-of-the-art GAN-based face synthesis PGGAN [7] where it shows reasonable classification performance while only using low dimensional features and a light model with fewer parameters.

2 RELATED WORKS

2.1 GAN-based Face Synthesis Methods

Since the inaugural work of [5], GANs has revolutionized image synthesis methods. A GAN model is consisted of two neural networks, known as the generator and the encoder, that are trained in tandem. The generator takes random noises as input and synthesizes an image, which is sent to a discriminator, aiming to differentiate synthesized images from the real ones. The two networks are trained to compete with each other: the generator aims to create ever more realistic images to defeat the classifier while the discriminator network is trained to be more effective in differentiating the two types of images. The training ends when the two networks reach an equilibrium of the game. The original GAN model has since experienced many improvements. In particular, to improve the stability in training, Radford et al. optimized the network architecture by introducing the deep convolutional GANs (DCGAN)

¹Fukuwarai is a traditional game played in Japan during the new year time. A player of Fukuwarai is blindfolded and is requested to put parts of the face (*i.e.*, the eyes, eyebrows, nose and mouth), usually printed on paper, onto a blank face.



Figure 1: Over the years, GAN models have been improved significantly over the quality of faces they synthesize. Here we show a few examples of different GAN models (a) GAN [5], (b) DCGAN [15], (c) COGAN [11], (d) PGGAN [7], (e) Style-GAN [8].

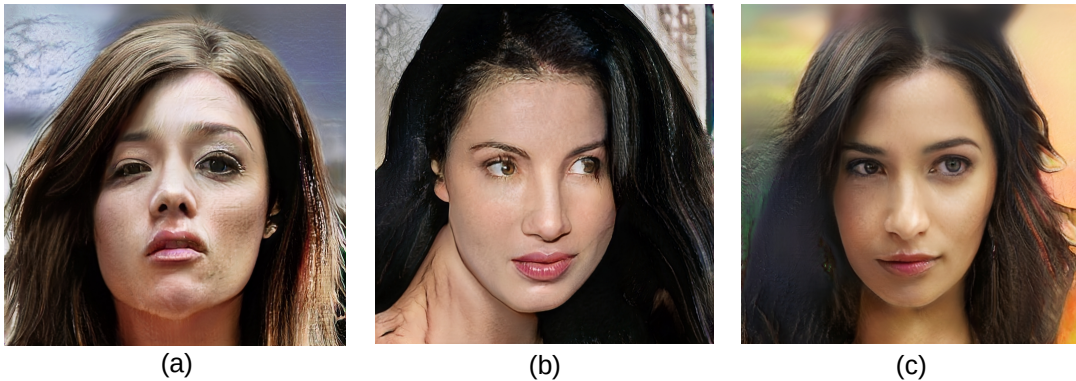


Figure 2: Abnormalities of GAN-based faces synthesized by PGGAN [7]. (a) In-symmetric size and location of eyes. (b) Mouth shifts towards left relatively to the nose. (c) Sharp and downwards lateral canthus (inner corner of eye) on left eye.

[15]. Coupled Generative Adversarial Networks (COGAN) learnt joint distribution from different domains further improved realism of the synthesized images [11]. However, the instability in the training process remains [2, 6, 9, 17, 18], which propagates to the synthesized samples and limits the model to only synthesize low resolution, see Figure 1.

PGGAN [7] is a major breakthrough for synthesizing high resolution realistic face images. It grew both generator and discriminator progressively, starting with generating 4×4 resolution images from generator. This generated image along with the training image resized into the same scale is feed into the discriminator. After the network are stabilized, a three layer blocks (similar to residual blocks), generating images with doubled heights and widths, faded into the network. These model stabilization through training and higher resolution layers fading in was carried out alternatively, until 1024×1024 resolution of the generated images is achieved. This approach not only improved training speed and stability, but also synthesized high resolution face images (1024×1024) with unprecedented fine details. The PGGAN model is further improved by style-transfer GAN (STGAN) [8], which treats face synthesis

problem as transferring styles of one face to another. However, STGAN is fundamentally different from previous GAN-based image synthesis models in that images of the best quality are generated conditioned on existing images instead of directly from random noises. Because of this reason, we do not consider detecting STGAN generated images in this work.

2.2 Detection Methods for GAN-synthesized Images

Compared to popularity of exploring strategies for synthesizing face images with GANs, methodologies to differentiate the real and synthesized images are far from satisfactory. Li et al [10] observed the color mismatch in H, S, V and Cb, Cr, Y channels between real and GAN-generated images. Similarly, McCloskey and Albright identified the frequency of saturated pixels and color image statistics of the GAN-generated images are different from the ones captured by cameras [13]. However, this color disparity could easily be removed by post processing after the image synthesis. On the other hand, Mo et al [14] and Tariq [19] designed deep convolutional neural networks classifiers for fake exposure, which usually

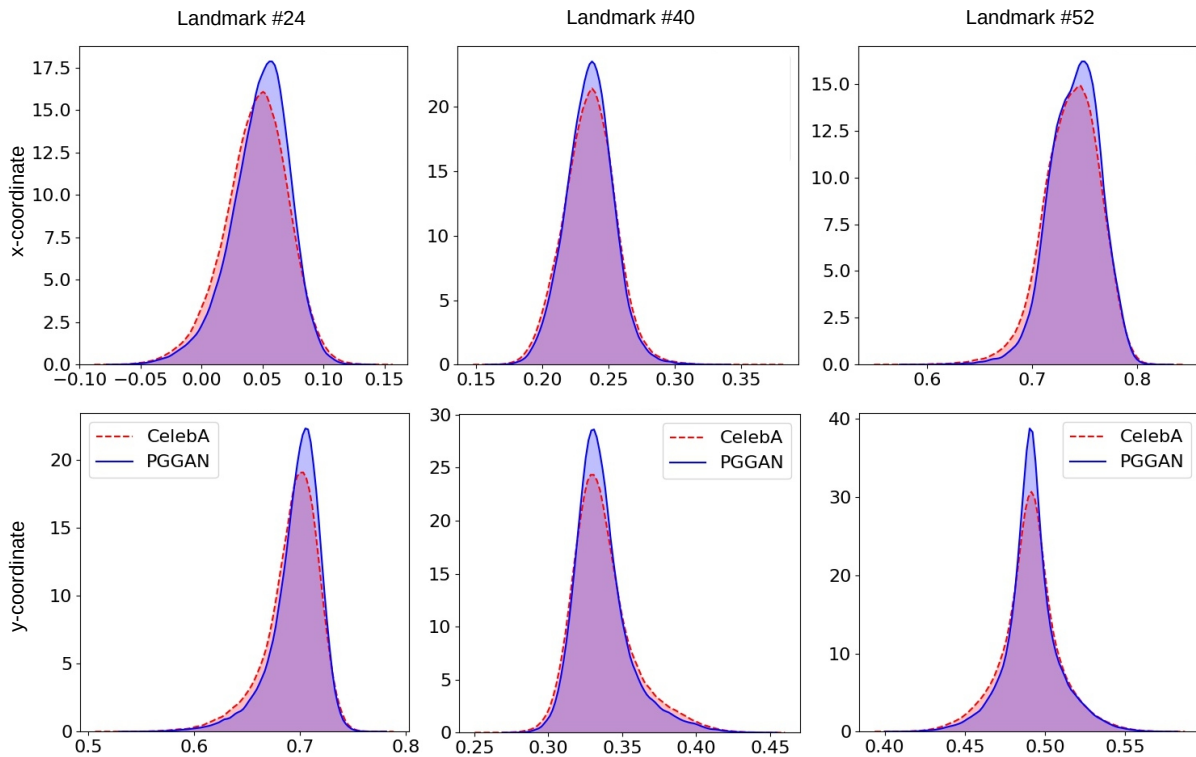


Figure 3: Density distribution of normalized face landmark locations on real (CelebA) and GAN-synthesized fake (PGGAN) faces over x - and y -coordinates. Real ones are from CelebA dataset with 200k+ images, and fake faces are from PGGAN dataset with 100k images.

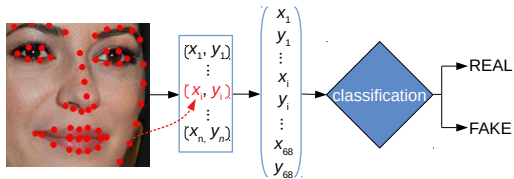


Figure 4: Pipeline for image our classification method. 68 landmarks are detected from face images which is warped into a standard configuration, followed by flattening landmark locations into 136D vector for classification.

requires GPU in training and testing, and not be able to reveal the mechanisms behind the classification.

3 METHOD

As we described in the Introduction, GAN-synthesized faces may exhibit inconsistent configurations of facial parts due to the weak global constraints. Several examples of this phenomenon are shown in Figure 2 for high resolution face images synthesized with the state-of-the-art PGGAN method [7]. In (a), we observe that the synthesized two eyes, nose and upper lips are not symmetric. In (b), the right eye is distorted and the mouth is shifted left-ward with regards to the tip of the nose. In (c), the face shows an unnatural lateral canthus (sharp and down ward inner corner of left eye) and different sizes of two eyes.

To quantify such inconsistencies, we compare facial landmark locations detected over GAN-synthesized and real faces. We first

run a face detector and extract facial landmarks, Figure 4. The detected landmarks are warped into a standard configuration in the region of $[0, 1] \times [0, 1]$ through an affine transformation by minimizing the alignment errors. To reduce the effect of face shape to the alignment result, we follow the standard procedure to estimate the warping transform using only facial landmarks in the central area of the face excluding those on the face contour. Figure 3 shows the differences in the aligned landmark locations for the real and GAN-synthesized faces in terms of their distributions along the x - and y - image coordinates. As these results show, the marginal distributions of landmarks for the GAN generated faces exhibit some consistent differences, and such differences are more prominent when we consider the joint distribution over all the coordinates of the ensemble of face landmark points. Therefore, we can use the vector formed by vectorizing all these landmark locations as a feature vector to build a classification system for differentiating GAN-synthesized and real faces, Figure 4.

There are three advantages of using such features for the classification tasks. First, this feature has relatively low dimension (it is twice the number of landmarks we extract from each face). This facilitates the construction of simpler classification schemes. Second, landmark locations are indifferent of image sizes, so there is no need to rescale the image in training and using the obtained classification method, which may also avoid the undesirable side-effect that leads the classifier to capture the artificial differences in image resolution due to the resizing operation. Third, the abnormalities of facial landmark locations attribute to the underlying

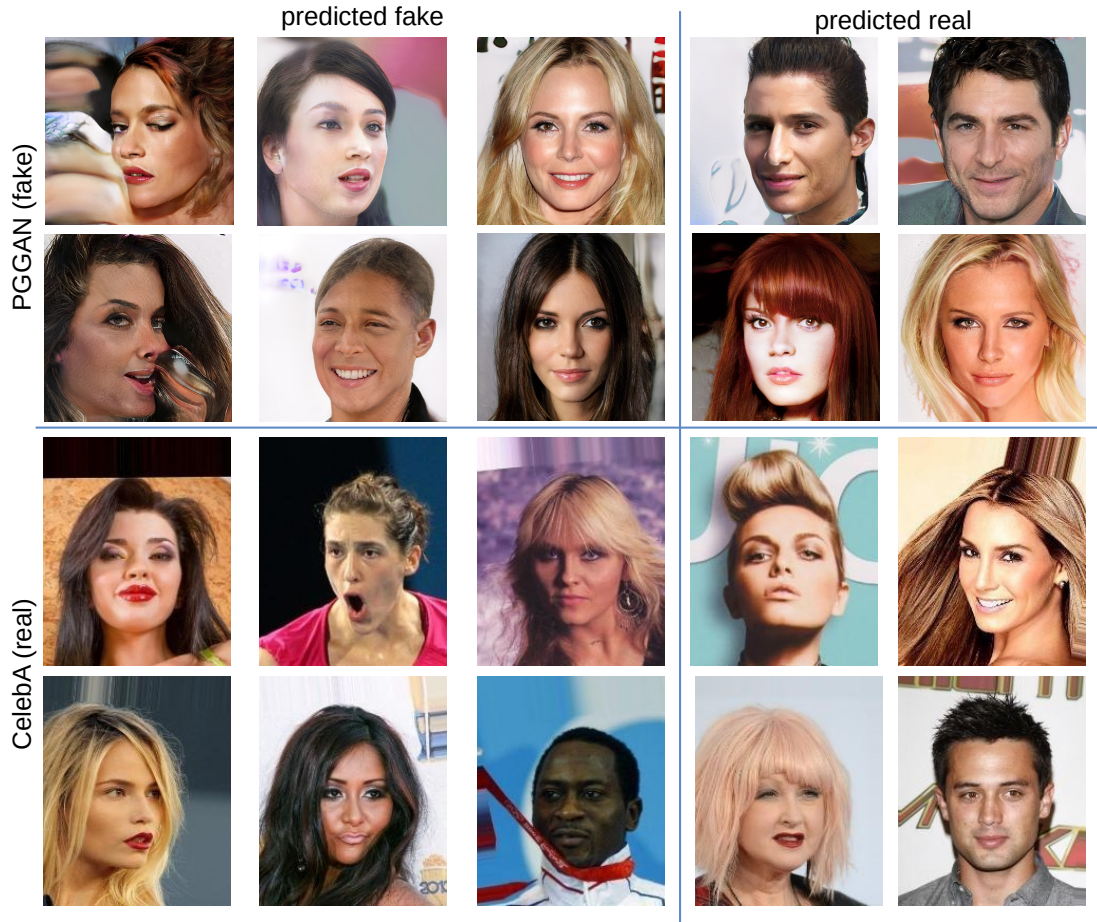


Figure 5: Examples of correct and incorrect predictions on CelebA and PGGAN datasets.

fundamental mechanism of GAN image synthesis, so it may not be trivially fixable without introducing more complex constraints into the GAN framework.

4 EXPERIMENTS

In this section, we report the experimental evaluations using landmark locations as features to distinguish real face images from the ones synthesized by GAN.

Choosing GAN-based Face Synthesis Method. Although there are a plethora of GAN-based face synthesis methods [5, 7, 8, 11, 15], we choose in this work the recent PGGAN to construct and evaluate the landmark location based classification method for classifying current state-of-the-art high quality GAN-synthesized faces. This choice is motivated by the following reasons. Early GAN-based face synthesis methods [5, 11, 15] produce low quality face images with low resolutions, so they are not representative to the state of the art. On the other hand, the most recent STGAN does not synthesize face images from noise as all other GAN-based method but treat it as a style-transfer problem.

Dataset. The training and testing of the SVM classifiers for exposing GAN-synthesized images are based on two datasets: (a)

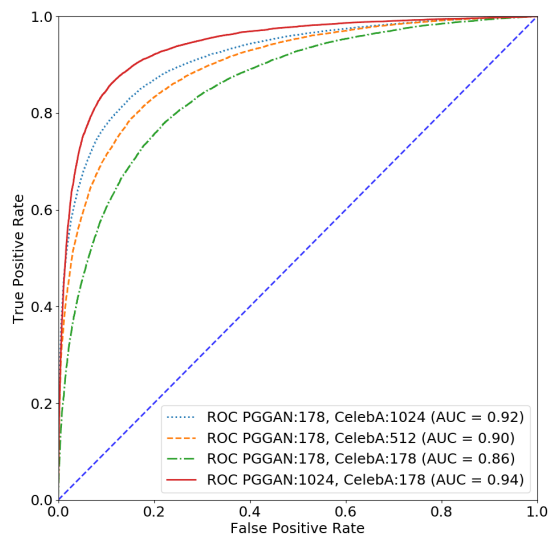
CelebFaces Attributes Dataset (CelebA) [12] contains more than 200K real face images with fixed resolution of 216×178 pixels. (b) PGGAN dataset [7] consisting of 100K PGGAN-synthesized face images at a resolution of 1024×1024 pixels are used as fake faces. 75% of both datasets are merged as negative and positive samples for training, and the rest 25% are used for testing.

Preprocessing and Training. Using the normalized locations of all face landmarks as features, we can develop a simple classification scheme to differentiate the real and GAN-synthesized faces, with standard classification methods such as SVM or neural networks, Figure 4. In this study, the normalized landmark locations of each face ($\in R^{68 \times 2}$) are flattened in to a vector ($\in R^{136 \times 1}$), which is standardized by subtracting the mean and divided by the standard deviation of all training samples. We trained SVM classifiers with radial basis function (RBF) kernel with a grid search on the hyperparameters using 5-fold cross validation. The losses of two classes are balanced by adjusting the sample loss inversely proportional to class frequencies in training dataset.

Performance. Figure 5 shows some examples of prediction results on PGGAN and CelebA datasets. PGGAN-synthesized faces with artifacts could be correctly predicted as fake faces, and the ones falsely predicted to be real mostly bears no visible defects. For

Table 1: AUCROC performance of our method and other deep neural network methods [19] on PGGAN and CelebA resized into different image widths.

method	# parameters	dataset (image resolution)		AUROC (%)
		CelebA (216 × 178)	PGGAN (1024 × 1024)	
VGG19	~143.7M	resize to (1243 × 1024)	remain at (1024 × 1024)	60.13
XceptionNet	~22.9M			85.03
NASNet	~3.3M			96.55
ShallowNetV3	-			99.99
Our method (SVM)	~110K			91.21
Our method (SVM)	~110K	original size (216 × 178)	original size (1024 × 1024)	94.13

**Figure 6: Performance by varying the widths of the PGGAN and CelebA images. Widths of resized images for each dataset are in the legend.**

the real face images in the CelebA dataset, some faces are falsely predicted to be fake. This may result from difficulties in accurately estimating landmark locations in faces with strong facial expression and occlusion, as shown in the bottom left figures in Figure 5.

Quantitative results of our method is shown in Table 1, in terms of the Area Under ROC (AUROC). As a comparison, we also include performance with different neural network architectures from [19] on the same dataset². Note that all methods in [19] take the image as input. To accommodate the different sizes of input images between the CelebA (216 × 178) and PGGAN (1024 × 1024) datasets, the images are resized to the same size and the results on enlarging the celebA images to 1243 × 1024 are reported in Table 1.

As the results show, the SVM classifier achieves an AUROC of 94.13% and outperforms several deep neural network based methods (e.g., VGG19 and XceptionNet). The two deep neural network based methods achieving higher classification accuracy are with much higher number of parameters. More importantly, these results

²These results are taken from the published paper [19] directly, because no code is currently available. The training and testing data may differ.

are obtained on resized images and no study was conducted on the effect of the resizing on the final classification – as upsampling an image lead to certain artifacts. It is not clear how much of the high performance can be attribute to the intrinsic difference between the two types of images. As we mentioned previously, the feature based on locations of facial landmarks is independent from image sizes and we compare the effect of resizing the two classes of images in another set of experiments shown in Figure 6, which demonstrates that the classification performance is relatively indifferent to the resizing operation. We would also like to emphasize that all these CNN models requires GPU for training and testing, while our method has much fewer parameters and only CPU for training and testing.

5 RESULTS ON FACE FORENSICS

Although the feature we proposed is originally designed for GAN-based face synthesis, we believe that other types of face synthesis methods may also exhibit similar abnormalities. To this end, we test our method on the FaceForensics dataset, which contains pairs of real and falsified videos synthesized by Face2Face [20]. It has 740 pairs of videos (726,270 images) for training, 150 pairs (151,052 images) for validation, and 150 videos (155,490 images) for testing [16]. The video frames in FaceForensics dataset vary in 576 to 1920 pixels for width and 480 to 1080 pixels for height. SVM models based on the landmark features are trained similarly on this dataset and we report the performance in Figure 7, which achieves a 0.83 AUROC for classification in individual frames. By averaging the classification prediction on individual videos, the AUROC increases to 0.90.

6 CONCLUSION

In this work, we proposed using aligned facial landmark locations as features to distinguish PGGAN synthesized fake human face images. Our method is based on the observation that current GAN-based algorithms uses random noises as input, which is good at depicting the details of face parts, but lack of constrains on the configuration of different face components. Consequently, it introduces errors in facial parts locations, which is non-trivial to be fixed in GAN models. We performed experiments to demonstrate this phenomenon and further developed classification models on this cue. The results indicated the effectiveness of our methods with low dimensional input, light-weighted models, and robust to scale variation.

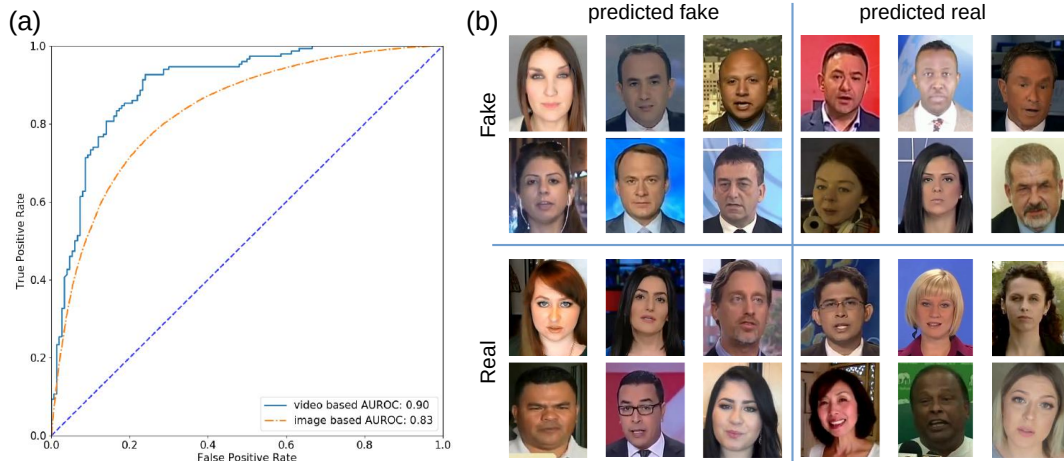


Figure 7: Classification on FaceForensics Testing Dataset. (a) The ROC curve and AUROC scores. (b) Examples of correct and incorrect predictions for both classes.

REFERENCES

- [1] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. 2017. A brief survey of deep reinforcement learning. *arXiv preprint arXiv:1708.05866* (2017).
- [2] David Berthelot, Thomas Schumm, and Luke Metz. 2017. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717* (2017).
- [3] Li Deng. 2014. A tutorial survey of architectures, algorithms, and applications for deep learning. *APSIPA Transactions on Signal and Information Processing* 3 (2014).
- [4] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NIPS*.
- [6] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. 2017. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*. 5767–5777.
- [7] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2017. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196* (2017).
- [8] Tero Karras, Samuli Laine, and Timo Aila. 2018. A style-based generator architecture for generative adversarial networks. *arXiv preprint arXiv:1812.04948* (2018).
- [9] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. 2017. How to train your DRAGAN. *arXiv preprint arXiv:1705.07215* 2, 4 (2017).
- [10] Haodong Li, Bin Li, Shunquan Tan, and Jiwu Huang. 2018. Detection of deep network generated images using disparities in color components. *arXiv preprint arXiv:1808.07276* (2018).
- [11] Ming-Yu Liu and Oncel Tuzel. 2016. Coupled generative adversarial networks. In *Advances in neural information processing systems*. 469–477.
- [12] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- [13] Scott McCloskey and Michael Albright. 2018. Detecting GAN-generated Imagery using Color Cues. *arXiv preprint arXiv:1812.08247* (2018).
- [14] Huaxiao Mo, Bolin Chen, and Weiqi Luo. 2018. Fake Faces Identification via Convolutional Neural Network. In *Proceedings of the 6th ACM Workshop on Information Hiding and Multimedia Security*. ACM, 43–47.
- [15] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434* (2015).
- [16] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. 2018. Faceforensics: A large-scale video dataset for forgery detection in human faces. *arXiv preprint arXiv:1803.09179* (2018).
- [17] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. In *Advances in neural information processing systems*. 2234–2242.
- [18] Tim Salimans and Durk P Kingma. 2016. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*. 901–909.
- [19] Shahroz Tariq, Sangyup Lee, Hoyoung Kim, Youjin Shin, and Simon S Woo. 2018. Detecting both machine and human created fake face images in the wild. In *Proceedings of the 2nd International Workshop on Multimedia Privacy and Security*. ACM, 81–87.
- [20] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2387–2395.
- [21] Qingchen Zhang, Laurence T Yang, Zhikui Chen, and Peng Li. 2018. A survey on deep learning for big data. *Information Fusion* 42 (2018), 146–157.