

Detection, Attribution and Localization of GAN Generated Images

Michael Goebel
mgoebel@ucsb.edu
University of California, Santa
Barbara
Santa Barbara, California

Lakshmanan Nataraj
nataraj@mayachitra.com
Mayachitra, Inc.
Santa Barbara, California

Tejaswi Nanjundaswamy
tejaswi@mayachitra.com
Mayachitra, Inc.
Santa Barbara, California

Tajuddin Manhar Mohammed
mohammed@mayachitra.com
Mayachitra, Inc.
Santa Barbara, California

Shivkumar Chandrasekaran
shiv@mayachitra.com
Mayachitra, Inc.
University of California
Santa Barbara, California

B.S. Manjunath
manj@mayachitra.com
Mayachitra, Inc.
University of California
Santa Barbara, California

ABSTRACT

Recent advances in Generative Adversarial Networks (GANs) have led to the creation of realistic-looking digital images that pose a major challenge to their detection by humans or computers. GANs are used in a wide range of tasks, from modifying small attributes of an image (StarGAN [14]), transferring attributes between image pairs (CycleGAN [91]), as well as generating entirely new images (ProGAN [36], StyleGAN [37], SPADE/GauGAN [64]). In this paper, we propose a novel approach to detect, attribute and localize GAN generated images that combines image features with deep learning methods. For every image, co-occurrence matrices are computed on neighborhood pixels of RGB channels in different directions (horizontal, vertical and diagonal). A deep learning network is then trained on these features to detect, attribute and localize these GAN generated/manipulated images. A large scale evaluation of our approach on 5 GAN datasets comprising over 2.76 million images (ProGAN, StarGAN, CycleGAN, StyleGAN and SPADE/GauGAN) shows promising results in detecting GAN generated images.

KEYWORDS

Image Forensics, Media Forensics, GAN Image Detection, GAN Image Localization, GAN Image Attribution, Detection of Computer Generated Images

1 INTRODUCTION

The advent of Convolutional Neural Networks (CNNs) [42, 71] has shown application in a wide variety of image processing tasks, and image manipulation is no exception. In particular, Generative Adversarial Networks (GANs) [24] have been one of the most promising advancements in image enhancement and manipulation - the generative Artificial Intelligence (AI) patents grew by 500% in 2019 [2]. Due to the success of using GANs for image editing, it is now possible to use a combination of GANs and off-the-shelf image-editing tools to modify digital images to such an extent that it has become difficult to distinguish doctored images from normal ones. In December 2019, Facebook announced that it removed

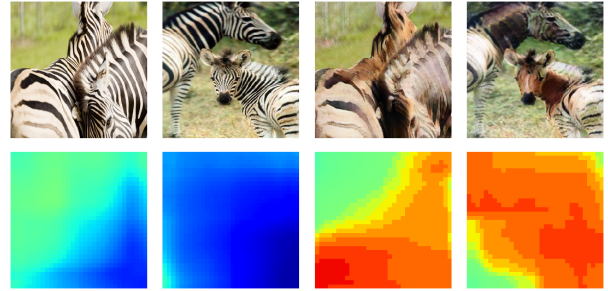


Figure 1: Input test set images on the top row, and our proposed detection heatmaps on the bottom. The two images on the left are authentic zebra images, those on the right are generated using CycleGAN.

hundreds of accounts whose profile pictures were generated using AI [1, 3].

The GAN training procedure involves a generator and discriminator. The generator may take in an input image and a desired attribute to change, then output an image containing that attribute. The discriminator will then try to differentiate between images produced by the generator and the authentic training examples. The generator and discriminator are trained in an alternate fashion, each attempting to optimize its performance against the other. Ideally, the generator will converge to a point where the output images are so similar to the ground truth that a human will not be able to distinguish the two. In this way, GANs have been used to produce “fake” images that are very close to the real input images. These include image-to-image attribute transfer (CycleGAN [91]), generation of facial attributes and expressions (StarGAN [14]), as well as generation of whole new images such as faces (ProGAN [36], StyleGAN [37]), indoors (StyleGAN) and landscapes (SPADE/GauGAN [64]). In digital image forensics, the objective is to both detect these fake GAN generated images, localize areas in an image which have been generated by GANs, as well as identify which type of GAN was used in generating the fake image.

In the GAN training setup, the discriminator functions directly as a classifier of GAN and non-GAN images. So the question could

be raised as to *why not use the GAN discriminator to detect if it's real or fake?* To investigate this, we performed a quick test using the CycleGAN algorithm under the maps-to-satellite-images category, where fake maps are generated from real satellite images, and vice versa. In our test, we observed that the discriminator accuracy over the last 50 epochs was only 80.4%. However, state-of-the-art deep learning detectors for CycleGAN often achieve over 99% when tested on the same type of data which they are trained [55, 61, 88]. Though the discriminator fills its role of producing a good generator, it does not compare performance wise to other methods which have been suggested for detection.

While the visual results generated by GANs are promising, the GAN based techniques alter the statistics of pixels in the images that they generate. Hence, methods that look for deviations from natural image statistics could be effective in detecting GAN generated fake images. These methods have been well studied in the field of steganalysis which aims to detect the presence of hidden data in digital images. One such method is based on analyzing co-occurrences of pixels by computing a co-occurrence matrix. Traditionally, this method uses hand crafted features computed on the co-occurrence matrix and a machine learning classifier such as support vector machines determines if a message is hidden in the image [72, 73]. Other techniques involve calculating image residuals or passing the image through different filters before computing the co-occurrence matrix [17, 23, 66]. Inspired by steganalysis and natural image statistics, we propose a novel method to identify GAN generated images using a combination of pixel co-occurrence matrices and deep learning. Here we pass the co-occurrence matrices directly through a deep learning framework and allow the network to learn important features of the co-occurrence matrices. This also makes it difficult to perform adversarial perturbations on the co-occurrence matrices since the underlying statistics will be altered. We also avoid computation of residuals or passing an image through various filters which results in loss of information. We rather compute the co-occurrence matrices on the image pixels itself. For detection, we consider a two class framework - real and GAN, where a network is trained on co-occurrence matrices computed on the whole image to detect if an image is real or GAN generated. For attribution, the same network is trained in a multi-class setting depending on which GAN the image was generated from. For localization, a network is trained on co-occurrence matrices computed on image patches and a heatmap was generated to indicate which patches are GAN generated. Detailed experimental results on large scale GAN datasets comprising over 2.76 million images originating from multiple diverse and challenging datasets generated using GAN based methods show that our approach is promising and will be an effective method for tackling future challenges of GANs.

The main contributions of the paper are as follows:

- We propose a new method for detection, attribution and localization of GAN images using a combination of deep learning and co-occurrence matrices.
- We compute co-occurrence matrices on different directions of an image and then train them using deep learning. For detection and attribution, the matrices are computed on the whole image and for localization, the matrices are computed on image patches to obtain a heatmap.

- We perform our tests on over 2.7 million images, which to our knowledge, is the largest evaluation on detection of GAN images.
- We provide explainability of our approach using t-SNE visualizations on different GAN datasets.
- We show the method holds under both varying JPEG compression factors and image patch sizes, accommodating a range of real-world use cases.

2 RELATED WORK

Since the seminal work on GANs [24], there have been several hundreds of papers on using GANs to generate images. These works focus on generating images of high perceptual quality [5, 25, 33, 36, 59, 67, 70], image-to-image translations [33, 85, 91], domain transfer [40, 76], super-resolution [43], image synthesis and completion [32, 46, 83], and generation of facial attributes and expressions [14, 40, 48, 65]. Several methods have been proposed in the area of image forensics over the past years [9, 21, 47, 53, 80]. Recent approaches have focused on applying deep learning based methods to detect tampered images [6–8, 11, 17, 69, 90].

In digital image forensics, detection of GAN generated images has been an active topic in recent times and several papers have been published in the last few years [4, 10, 10, 13, 20, 22, 26, 29, 30, 34, 35, 39, 44, 45, 55–58, 60, 61, 63, 77, 79, 81, 82, 86–88, 92]. Other similar research include detection of computer generated (CG) images [19, 52, 68, 84]

In [55], Marra et al. compare various methods to identify CycleGAN images from normal ones. The top results they obtained are using a combination of residual features [16, 17] and deep learning [15]. In [45], Li et al. compute the residuals of high pass filtered images and then extract co-occurrence matrices on these residuals, which are then concatenated to form a feature vector that can distinguish real from fake GAN images. In [88], Zhang et al. identify an artifact caused by the up-sampling component included in the common GAN pipeline and show that such artifacts are manifested as replications of spectra in the frequency domain and thus propose a classifier model based on the spectrum input, rather than the pixel input.

We had previously proposed a 3 channel co-occurrence matrix based method [61], and many other papers have shown the efficacy of this method in their experimental evaluations [31, 50, 54, 62, 63, 78, 87]. However, in this paper we compute co-occurrence matrices on horizontal, vertical and diagonal directions, as well as compute them on image patches, thus facilitating detection, attribution and localization of GAN generated images.

3 METHODOLOGY

3.1 Co-Occurrence Matrix Computation

The co-occurrence matrices represent a two-dimensional histogram of pixel pair values in a region of interest. The vertical axis of the histogram represents the first value of the pair, and the horizontal axis, the second value. Equation 1 shows an example of this computation for a vertical pair.

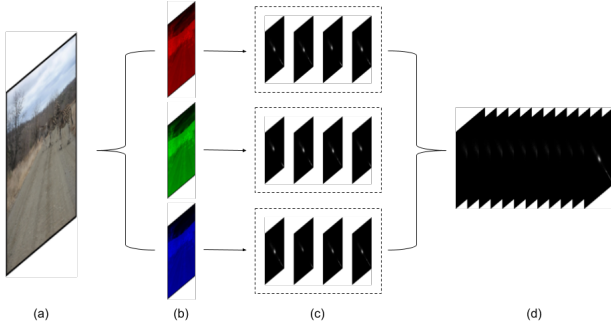


Figure 2: An example co-occurrence computation. The input image (a) is split into its three color channels (b). For each color channel, 4 different pairs of pixels are used to generate 2-dimensional histograms (c). Horizontal, vertical, diagonal, and anti-diagonal pairs are considered. These histograms are then stacked to produce a single tensor (d). For some tests, only a subset of the co-occurrence matrices will be used.

$$C_{i,j} = \sum_{m,n} \begin{cases} 1, & I[m,n] = i \text{ and } I[m+1,n] = j \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Under the assumption of 8-bit pixel depth, this will always produce a co-occurrence matrix of size 256×256 . This is a key advantage of such a method, as it will allow for the same network to be trained and tested on a variety of images without resizing.

Which pairs of pixels to take was one parameter of interest in our tests. For any pixel not touching an edge, there are 8 possible neighbors. We consider only 4 of these for our tests; right, bottom right, bottom, and bottom left. The other 4 possible pairs will provide redundant information. For example, the left pairs are equivalent to swapping the order of the first and second pixel in the right pair. In the co-occurrence matrix, this corresponds to a simple transpose. There are many subsets of these 4 pairs which could be taken, but our tests consider only a few; horizontal, vertical, horizontal and vertical, or all.

Before passing these matrices through a CNN, some pre-processing is done. First, each co-occurrence matrix is divided by its maximum value. Given that the input images may be of varying sizes, this will force all inputs into a consistent scale. After normalization, all co-occurrence matrices for an image are stacked in the depth dimension. In the example of an RGB image with all 4 co-occurrence pairs, this will produce a new image-like feature tensor of size $256 \times 256 \times 12$. Figure 2 gives a visualization of this process.

3.2 Convolutional Neural Networks

While the co-occurrence matrices are not themselves images, treating them as so has some theoretical backing. One of the primary motivations for using CNNs in image processing is their translation invariance property. In the case of a co-occurrence matrix, a translation along the main diagonal corresponds to adding a constant

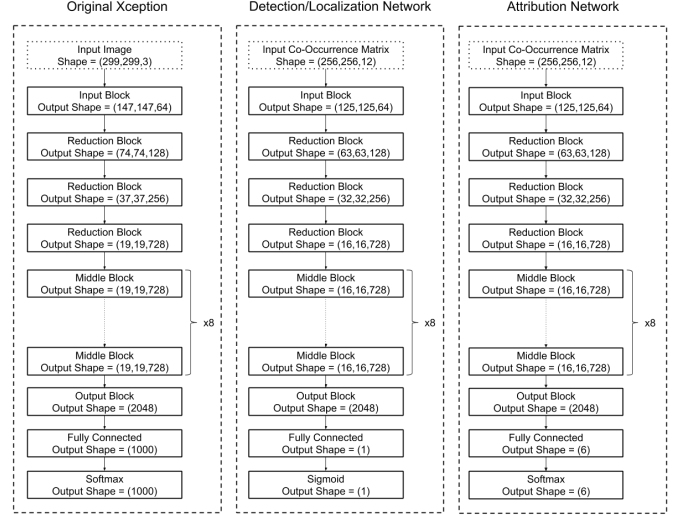


Figure 3: The original Xception network [15], shown next to our two modified models. Our architectures for detection and attribution are the same, except for the last layer and activation.

value to the image. We would not expect this manipulation to affect the forensic properties.

In this paper, we use Xception Net [15] deep neural network architecture for detection, attribution and localization of GAN generated images. The Xception network is a modified version of Inception network [74] but was created under a stronger theoretical assumption than the original Inception, where cross-channel correlations are completely split from spatial correlations by use depth-wise separable convolutions. The network also includes residual connections, as shown in Figure 3. For these reasons, the authors claim that Xception can more easily find a better convergence point than most other CNN architectures, while keeping model capacity low [15]. In this paper, we modify the original input and output shapes in the Xception network to accommodate our task as shown in Figure 3. The initial convolutional portions of the network remain unchanged, though the output sizes of each block are slightly different. This small change in size is accommodated by the global pooling step. Finally, the last fully connected layer of each network is changed to the desired number of output classes, and given the appropriate activation. For detection and attribution, our architectures are the same except for the last layer and activation. For localization, no changes were made to the model architecture but co-occurrence matrices were extracted on small image patches, and individually passed through the network.

4 DATASETS

We evaluated our method on five different GAN architectures, of which each was trained on several different image generation tasks: ProGAN [36], StarGAN [14], CycleGAN [91], StyleGAN [37], and SPADE/GauGAN [64]. The modifications included image-to-image

translation, facial attribute modification, style transfer, and pixel-wise semantic label to image generation. A summary of the datasets, including the number of images from each class, is shown in Figure 5. These comprise a total of more than 2.76 million images of which 1.69 million images are real images and 1.07 million images are fake GAN generated images. In several cases, one or more images in the GAN generated category will be directly associated with an image in the authentic class. For example, a person’s headshot untampered, blond, aged, and gender reversed will all be in the dataset. However, the splitting for training accounts for this, and will keep all of these images together to be put into either training, validation, or test. Some sample images from all the GAN datasets are shown in Figure 4.

4.0.1 StarGAN. This dataset consists of only celebrity photographs from the CelebA dataset [49], and their GAN generated counterparts [14]. The GAN changes attributes of the person to give them black hair, brown hair, blond hair, different gender, different age, different hair and gender, different hair and age, different gender and age, or different hair, age, and gender. These are the smallest of all of the training images, being a square of size 128 pixels.

4.0.2 CycleGAN. This datasets includes image-to-image translations between a wide array of image classes [91]. The sets horse2zebra, apple2orange, and summer2winter do a strict image-to-image translation, with the assumption that the GAN will learn the areas to modify. While the whole output is generated by the GAN, the changes for these will ideally be more localized. Ukiyoe, Vangogh, Cezanne, and Monet are four artists which the GAN attempts to learn a translation from photographs to their respective styles of painting. Facades and cityscapes represent the reverse of the image segmentation task. Given a segmentation map as input, they produce an image of a facade or cityscape. Map2sat takes in a Google Maps image containing road, building, and water outlines, and generates a hypothetical satellite image.

4.0.3 ProGAN. This dataset consists of images of celebrities, and their GAN generated counterparts, at a square size of 1024 pixels [36]. All data was obtained per the instructions provided in the paper’s Github repository.

4.0.4 SPADE/GauGAN. SPADE/GauGAN contains realistic natural images generated using GANs [64]. This dataset uses images from ADE20k [89] dataset containing natural scenes and COCO-Stuff [12] dataset comprising day-to-day images of things and other stuff, along with their associated segmentation maps. These untampered images are considered as real images in the GAN framework, and the pretrained models provided by the SPADE/GauGAN authors are used to generate GAN images from the segmentation maps.

4.0.5 StyleGAN. This dataset contains realistic images of persons, cars, cats and indoor scenes [37]. Images for this dataset were provided by the authors.

Table 1: Comparison of different popular ImageNet [18] classification architectures on classifying GANs from co-occurrence matrices. All datasets are mixed for training, validation, and testing. The features are extracted from a whole image, with no JPEG compression.

Network	Accuracy
VGG16 [71]	0.6115
ResNet50 [27]	0.9677
ResNet101 [27]	0.9755
ResNet152V2 [28]	0.9795
ResNet50V2 [28]	0.9856
InceptionResNetV2 [74]	0.9885
InceptionV3 [75]	0.9894
ResNet101V2 [28]	0.9900
Xception [15]	0.9916

5 EXPERIMENTS

5.1 Training Procedure

All deep learning experiments in this paper were done using Keras 2.2.5 and all training was done using an Adam optimizer [41], learning rate of 10^{-4} , and cross-entropy loss. A batch size of 64 was used for all experiments. Unless otherwise stated, a split of 90% training, 5% validation, and 5% test was used. Given the large amount of data available, a single iteration through the entire dataset for training took 10 hours on a single Titan RTX GPU. To allow for more frequent evaluation on the validation set, the length of an epoch was capped at 100 batches. Validation steps were also capped at 50 batches, and test sets at 2000 batches. After training for a sufficient period of time for the network to converge, the checkpoint which scored the highest in validation was chosen for testing. For experiments to determine hyper-parameters, training was capped at 50 epochs, and took approximately 3 hours each on a single Titan RTX. After determination of hyper-parameters, training of the final model was done for 200 epochs, taking approximately 12 hours.

5.2 Comparison with other CNN architectures:

First we evaluate our method on different well known CNN architectures: VGG16 [71], ResNet50 and ResNet101 [27], ResNet50V2, ResNet101V2 and ResNet152V2 [28], InceptionV3 and Inception-ResNetV2 [74], and Xception [15]. Shown in Table 1 are the results for the different CNN networks. Though designed for ImageNet classification, all models take in an image with height, width, and 3 channels, and output a one-hot encoded label. The models are used as-is, with the following slight modifications. First, the number of input channels is set to be the depth of the co-occurrence feature tensor. Second, input shape was fixed at 256x256. Third, the number of output channels was set to 1. All of these parameters were passed as arguments to the respective Keras call for each model. A small margin separated the top performers, though Xception was the best with an accuracy of **0.9916** and had fewer parameters than others. For this reason, we chose Xception for the remainder of the experiments.



Figure 4: Sample images from different GAN datasets (a) Real images and (b) GAN images from different GAN datasets (top to bottom): ProGAN [36], StarGAN [14], CycleGAN [91], StyleGAN [37], and SPADE/GauGAN [64].

	real	fake
base	1696998	1073582
└─ stargan	<u>3279</u>	<u>29511</u>
└─ celeba	3279	29511
└─ cyclegan	<u>18151</u>	<u>18151</u>
└─ map2sat	1096	1096
└─ ukiyoe	1500	1500
└─ vangogh	1500	1500
└─ horse2zebra	2401	2401
└─ cezanne	1500	1500
└─ cityscapes	2975	2975
└─ apple2orange	2014	2014
└─ summer2winter	2193	2193
└─ monet	2572	2572
└─ facades	400	400
└─ progan	<u>30000</u>	<u>74000</u>
└─ celeba_hq	30000	74000
└─ spade	<u>145497</u>	<u>145497</u>
└─ ade20k	22210	22210
└─ coco_stuff	123287	123287
└─ stylegan	<u>1500071</u>	<u>806423</u>
└─ bedroom_lsun	500000	278790
└─ cat_lsun	500071	279867
└─ car_lsun	500000	247766

Figure 5: Quantitative summary of the GAN datasets used in our experiments.

5.3 Comparison of Co-occurrence Matrix Pairs

Next we perform tests with different co-occurrence pairs, shown in Table 2. These experiments included JPEG compression, randomly selected from quality factors of 75, 85, 90, and no compression. Interestingly, it seems that the addition of more co-occurrence pairs did not significantly improve performance. For the remainder of the test, all 4 co-occurrence pairs were used.

5.4 Effect of patch size

For applications, the two parameters of interest were JPEG compression and patch size. The results for different patch sizes are

Table 2: Test on difference co-occurrence pairs. These were done on the whole image, with the additional challenge of JPEG compression. The JPEG quality factor was randomly selected with equal probability from the set of 75, 85, 90, or no JPEG compression

Pairs	Accuracy
Horizontal	95.51
Vertical	95.56
Hor and Ver	95.17
Hor, Ver, and Diag	95.68

Table 3: Accuracy when trained on one patch size, and tested on another. Data for training and testing has been pre-processed using JPEG compression with quality factors randomly selected from 75, 85, 90 or none.

		Train		
		64	128	256
Test	64	0.7814	0.7555	0.6778
	128	0.8273	0.8336	0.8158
	256	0.8311	0.8546	0.8922

shown in Table 3. These results are from images JPEG compressed by a factor randomly selected from 75, 85, 90, and none. A model is trained for each of the possible patch sizes, and then each model is tested against features from each patch size. It should be noted that in cases where the input image is smaller than the requested patch size, the whole image is used. There is notable generalization between different patch sizes, in that the model trained on a patch size of 256 and tested on 128 achieves an accuracy within a few percentage points of a model trained and tested on 128. Thus we would expect our models to work with a variety of untested patch sizes within a reasonable range while only taking a minor performance drop.

5.5 Effect of JPEG compression

Now assuming a fixed patch size of 128, we varied the JPEG quality factors: 75,85,90 and no compression. The model was again trained only on one particular JPEG factor as shown in Table 4. As expected,

Table 4: Test accuracy when model is trained on images pre-processed with one JPEG quality factor, and tested on another.

		Train			
		75	85	90	None
Test	75	0.7738	0.7448	0.7101	0.6605
	85	0.8209	0.8593	0.8362	0.7209
	90	0.8310	0.8690	0.8756	0.7651
	None	0.9198	0.9386	0.9416	0.9702

Table 5: Train on all but one GAN, test on the held out images. Patch size of 128, no JPEG compression.

Test GAN	Accuracy
StarGAN	0.8490
CycleGAN	0.7411
ProGAN	0.6768
SPADE	0.9874
StyleGAN	0.8265

we see that performance increases with respect to quality factor. However, this table also shows that the model does not overfit to a particular quality factor, in that testing on a slightly better or worse quality factor gives a score not far from a model tuned to the particular test quality factor.

5.6 Generalization

To test the generalization between GANs, leave-one-out cross validation was used for each GAN architecture. One dataset of GAN images is used for testing and remaining GAN image datasets are used for training. Here, a patch size of 128 was used with no JPEG compression. From Table 5, we see that some GAN datasets such as SPADE, StarGAN and StyleGAN have high accuracy and are more generalizable. However, the accuracies for CycleGAN and ProGAN are lower in comparison, thus suggesting that images from these GAN categories should not be discarded when building a bigger GAN detection framework. We also considered computing co-occurrence matrices on the whole image and then repeated the above experiment, but the overall accuracy did not improve.

Visualization using t-SNE: To further investigate the variability in the GAN detection accuracies under the leave-one-out setting, we use t-SNE visualization [51] from outputs of the penultimate layer of the CNN, using images from the test set (as shown in Figure 6). The t-SNE algorithm aims to reduce dimensionality of a set of vectors while preserving relative distances as closely as possible. With an L2 distance function and linear transformation, this can be efficiently found by PCA. While there are many solutions to this problem for different distance metrics and optimization methods, KL divergence on the Student-t distribution used in t-SNE has shown the most promising results on real-world data [51].

To limit computation time, no more than 1000 images were used for a particular GAN from either the authentic or GAN classes. As recommended in the original t-SNE publication, the vector was first reduced using Principle Component Analysis (PCA). The original

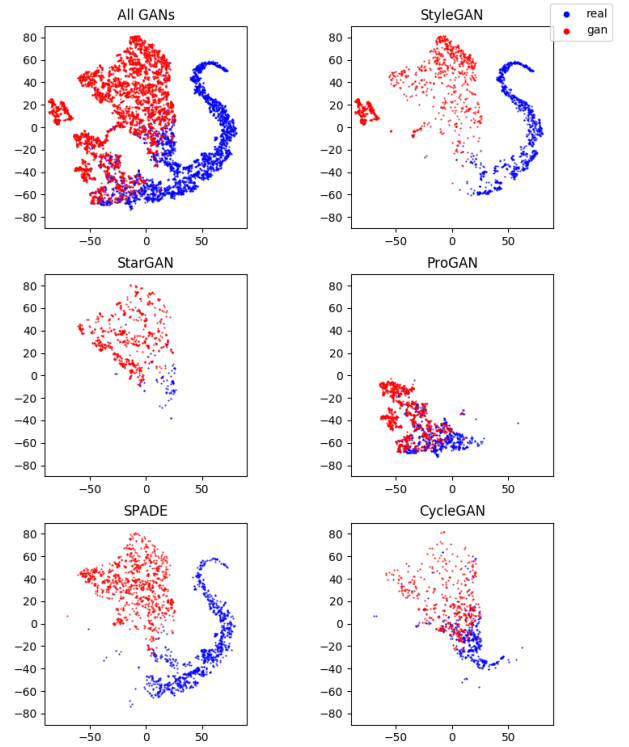


Figure 6: Visualization of images from different GAN datasets using t-SNE [51].

2048 were reduced to 50 using PCA, and passed to the t-SNE algorithm. As we see in Figure 6, the images in CycleGAN and ProGAN are more tightly clustered, thus making them difficult to distinguish between real and GAN generated images, while the images from StarGAN, SPADE and StyleGAN are more separable, thus resulting in higher accuracies in the leave-one-out experiment.

5.7 Comparison with State-of-the-art

We compare our proposed approach with various state-of-the-art methods [55, 61, 88] on the CycleGAN dataset. In [55], Marra et al. proposed the leave-one-category-out benchmark test to see how well their methods work when one category from the CycleGAN dataset is kept for testing and remaining are kept for training. The methods they evaluated are based on steganalysis, generic image manipulations, detection of computer graphics, a GAN discriminator used in the CycleGAN paper, and generic deep learning architecture pretrained on ImageNet [18], but fine tuned to the CycleGAN dataset. Among these the top performing ones were from steganalysis [16, 23] based on extracting features from high-pass residual images, a deep neural network designed to extract residual features [17] (Cozzolino2017) and XceptionNet [15] deep neural network trained on ImageNet but fine-tuned to this dataset. Apart from Marra et al. [55], we also compare our method with approaches including Nataraj et al. (Nataraj2019) [61], which uses co-occurrence matrices computed in the horizontal direction, and

Table 6: Comparison with State-of-the-art.

Method	ap2or	ho2zeb	wint2sum	citysc.	facades	map2sat	Ukiyoe	Van Gogh	Cezanne	Monet	Average
Steganalysis feat.	0.9893	0.9844	0.6623	1.0000	0.9738	0.8809	0.9793	0.9973	0.9983	0.9852	0.9440
Cozzalino2017	0.9990	0.9998	0.6122	0.9992	0.9725	0.9959	1.0000	0.9993	1.0000	0.9916	0.9507
XceptionNet	0.9591	0.9916	0.7674	1.0000	0.9856	0.7679	1.0000	0.9993	1.0000	0.9510	0.9449
Nataraj2019	0.9978	0.9975	0.9972	0.9200	0.8063	0.9751	0.9963	1.0000	0.9963	0.9916	0.9784
Zhang2019	0.9830	0.9840	0.9990	1.0000	1.0000	0.7860	0.9990	0.9750	0.9920	0.9970	0.9720
Proposed approach	0.9982	0.9979	0.9982	0.9366	0.9498	0.9776	0.9973	0.9980	0.9993	0.9697	0.9817

Table 7: Number of images per class

	Train	Val	Test
Authentic	1,612,202	42,382	42,397
StarGAN	28,062	738	711
CycleGAN	17,265	439	439
ProGAN	70,286	1833	1,881
SPADE	138,075	3,717	3,704
StyleGAN	766,045	20,220	20,158

Zhang et al.(Zhang2019) [88], which uses spectra of up-sampling artifacts used in the GAN generating procedure to classify GAN images.

Table 6 summarizes the results of our proposed approach against other state-of-the-art approaches. Our method obtained the best average accuracy of **0.9817**, when compared with other methods. Even on individual categories, our method obtained more than 0.90 on all categories.

5.8 Tackling newer challenges like StyleGAN2

Apart from generalization, we tested our method on 100,000 images from the recently released StyleGAN2 [38] dataset of celebrity faces. The quality of these images were much better than the previous version and appeared realistic. When we tested on this dataset without any fine-tuning, we obtained an accuracy of 0.9464. This shows that our approach is promising in adapting to newer challenges. We also fine-tuned to this dataset by adding 100,000 authentic images randomly chosen from different GAN datasets, thus our new dataset comprised of 100,000 authentic images and 100,000 StyleGAN2 images. Then, we split this data into 40% training, 10% validation and 50% testing. When we trained a new network on this dataset, we obtained a validation accuracy of 0.9984 and testing accuracy of 0.9972, thus also confirming that our approach can be made adjustable to newer GAN datasets.

5.9 GAN Attribution/Classification

While the primary area of interest is in determining the authenticity of an image, an immediate extension would be to determine which GAN was used. Here we perform an additional experiment on GAN class classification/attribution as a 6-class classification problem, the classes being: Real, StarGAN, CycleGAN, ProGAN, SPADE/GauGAN and StyleGAN. The number of output layers in the CNN was changed from 1 to 6, and output with the largest value was selected as the estimate. A breakdown of the number of images per class for training, validation and testing is given in Table 7. First, the network was trained where the input co-occurrence matrices were computed on the whole image. The training procedure was kept the same as with all other tests in the paper, with the exception of using a batch size of 60, and 10 images from each class per batch. This encouraged the network to not develop a bias

towards any particular GAN for which we have more training data. First we consider the images as they are provided in the datasets. The classification results are shown in the form of confusion matrices in Table 8. For convenience, we also report the equal prior accuracy, equal to the average along the diagonal of the confusion matrix. This equal prior accuracy can be interpreted as the classification accuracy if each class is equally likely. We obtain an overall classification accuracy (considering equal priors) of 0.9654. High classification accuracy was obtained for most categories. StyleGAN had comparatively lower accuracy but still more than 90%, being mostly confused with SPADE/GauGAN and CycleGAN. These results show that our approach can also be used to identify which category of GAN was used.

Next, we trained the network using a patch size of 128×128 as input, and repeated the experiment. This is to see how well our method can be used for detection, localization as well as classification. The classification results are shown in Table 9. Now, we obtain an overall classification accuracy (considering equal priors) of 0.8477 (a drop of 12% when compared to full image accuracy). High classification accuracy was obtained for StarGAN, CycleGAN and ProGAN, while SPADE/GauGAN and StyleGAN had comparatively lower accuracies. These could be due to many factors such as the number of test images per class, patch size, and the authentic image datasets that were used for training in generating these GAN images. In Table 10 we repeat the same experiment (with patch size 128×128) but with images that were randomly preprocessed with JPEG quality factors of 75, 85, 90, or no JPEG compression, with each of the four preprocessing methods equally likely. For this experiment, the overall classification accuracy drops slightly to 0.8088 due to the impact of JPEG compression.

For the multi-class experiment trained without JPEG compression, we repeat the t-SNE visualization procedure. Figure 7 shows all data-points on a single plot. These visualizations further support the results from the classification experiment.

5.10 Localization

Figure 8 show two example localization outputs. The image is processed in overlapping patches, with a particular stride and patch size. A co-occurrence matrix is then extracted for each patch, and passed through the CNN to produce a score. For pixels which are a part of multiple patches, the scores are simply the mean of all of the patch responses. These two examples use a patch size of 128, and a stride of 8. We can see that the heatmaps are predominantly blue for real images and predominantly red for GAN generated images. This further supports that our method can be effectively used for GAN localization.

Table 8: Confusion matrix on images from GAN datasets without any pre-processing on the full image. Equal prior accuracy of 0.9654.

		Predicted Label					
		Real	StarGAN	CycleGAN	ProGAN	SPADE	StyleGAN
GT Label	Real	0.975	0.000	0.000	0.016	0.002	0.006
	StarGAN	0.000	0.976	0.014	0.000	0.010	0.000
	CycleGAN	0.000	0.000	0.964	0.000	0.036	0.000
	ProGAN	0.000	0.000	0.000	1.000	0.000	0.000
	SPADE	0.001	0.000	0.019	0.000	0.975	0.005
	StyleGAN	0.007	0.000	0.022	0.000	0.068	0.902

Table 9: Confusion matrix on images from GAN datasets without any pre-processing on 128×128 patches. Equal prior accuracy of 0.8477.

		Predicted Label					
		Real	StarGAN	CycleGAN	ProGAN	SPADE	StyleGAN
GT Label	Real	0.826	0.003	0.016	0.021	0.066	0.068
	StarGAN	0.000	0.933	0.054	0.000	0.006	0.006
	CycleGAN	0.000	0.002	0.959	0.002	0.032	0.005
	ProGAN	0.000	0.002	0.008	0.981	0.004	0.005
	SPADE	0.001	0.025	0.210	0.008	0.728	0.029
	StyleGAN	0.003	0.025	0.101	0.009	0.203	0.659

Table 10: Confusion matrix with JPEG compression (128×128 patches). Equal prior accuracy of 0.8088. The images were preprocessed using a JPEG factor of 75, 85, 90, or no compression. Each of these four possible preprocessing functions was randomly selected with equal probability for every image.

		Predicted Label					
		Real	StarGAN	CycleGAN	ProGAN	SPADE	StyleGAN
GT Label	Real	0.741	0.005	0.020	0.026	0.103	0.104
	StarGAN	0.006	0.927	0.023	0.000	0.031	0.012
	CycleGAN	0.009	0.014	0.892	0.007	0.074	0.005
	ProGAN	0.002	0.003	0.009	0.973	0.007	0.007
	SPADE	0.075	0.015	0.095	0.009	0.765	0.042
	StyleGAN	0.114	0.021	0.059	0.008	0.243	0.555

6 CONCLUSIONS

In this paper, we proposed a novel method to detect and attribute GAN generated images, and localize the area of manipulations. Detailed experimental results using a collection of over 2.7 million GAN and authentic images encompassing 5 major GAN datasets demonstrate that the proposed model is highly effective on a range of image scales and JPEG compression factors. In addition, the t-SNE visualization with the neural network deep features showed promising separation of GAN and authentic images using our method.

ACKNOWLEDGMENTS

This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA). The views, opinions and/or findings expressed are those of the author and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government.

REFERENCES

- [1] [n.d.]. Facebook removes bogus accounts that used AI to create fake profile pictures. <https://www.cnet.com/news/facebook-removed-fake-accounts-that-used-ai-to-create-fake-profile-pictures/>. <https://www.cnet.com/news/facebook-removed-fake-accounts-that-used-ai-to-create-fake-profile-pictures/>
- [2] [n.d.]. Patent Filings for Generative AI Have Grown 500% This Year as Brands Test Its Potential. <https://www.adweek.com/digital/patent-filings-for-generative-ai-have-grown-500-this-year-as-brands-test-its-potential/>. <https://www.adweek.com/digital/patent-filings-for-generative-ai-have-grown-500-this-year-as-brands-test-its-potential/>
- [3] [n.d.]. Removing Coordinated Inauthentic Behavior From Georgia, Vietnam and the US. <https://about.fb.com/news/2019/12/removing-coordinated-inauthentic-behavior-from-georgia-vietnam-and-the-us/>. <https://about.fb.com/news/2019/12/removing-coordinated-inauthentic-behavior-from-georgia-vietnam-and-the-us/>
- [4] Michael Albright, Scott McCloskey, and ACST Honeywell. 2019. Source Generator Attribution via Inversion. *arXiv preprint arXiv:1905.02259* (2019).
- [5] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein gan. *arXiv preprint arXiv:1701.07875* (2017).
- [6] Jawadul H Bappy, Amit K Roy-Chowdhury, Jason Bunk, Lakshmanan Nataraj, and BS Manjunath. 2017. Exploiting Spatial Structure for Localizing Manipulated Image Regions. In *Proceedings of the IEEE International Conference on Computer Vision*.
- [7] Belhassen Bayar and Matthew C Stamm. 2016. A deep learning approach to universal image manipulation detection using a new convolutional layer. In *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security*. 5–10.

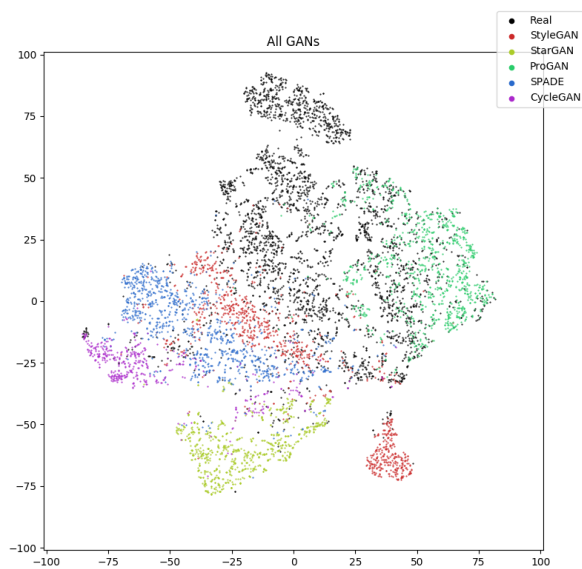


Figure 7: t-SNE visualization of 6 classes: Real, StyleGAN, StarGAN, ProGAN, SPADE/GauGAN and CycleGAN

[8] Belhassen Bayar and Matthew C Stamm. 2017. Design Principles of Convolutional Neural Networks for Multimedia Forensics. In *The 2017 IS&T International Symposium on Electronic Imaging: Media Watermarking, Security, and Forensics*. IS&T Electronic Imaging.

[9] Gajanan K Birajdar and Vijay H Mankar. 2013. Digital image forgery detection using passive techniques: A survey. *Digital Investigation* 10, 3 (2013), 226–245.

[10] Nicolò Bonettini, Paolo Bestagini, Simone Milani, and Stefano Tubaro. 2020. On the use of Benford’s law to detect GAN-generated images. *arXiv preprint arXiv:2004.07682* (2020).

[11] Jason Bunk, Jawadul H Bappy, Tajuddin Manhar Mohammed, Lakshmanan Nataraj, Arjuna Flenner, BS Manjunath, Shivkumar Chandrasekaran, Amit K Roy-Chowdhury, and Lawrence Peterson. 2017. Detection and Localization of Image Forgeries using Resampling Features and Deep Learning. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*. IEEE, 1881–1889.

[12] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. 2018. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1209–1218.

[13] Zehao Chen and Hua Yang. 2020. Manipulated Face Detector: Joint Spatial and Frequency Domain Attention Network. *arXiv preprint arXiv:2005.02958* (2020).

[14] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. 2018. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8789–8797.

[15] François Chollet. 2017. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint* (2017), 1610–02357.

[16] Davide Cozzolino, Diego Gragnaniello, and Luisa Verdoliva. 2014. Image forgery detection through residual-based local descriptors and block-matching. In *Image Processing (ICIP), 2014 IEEE International Conference on*. IEEE, 5297–5301.

[17] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. 2017. Recasting residual-based local descriptors as convolutional neural networks: an application to image forgery detection. In *Proceedings of the 5th ACM Workshop on Information Hiding and Multimedia Security*. ACM, 159–164.

[18] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*.

[19] Ahmet Emir Dirik, Sevinc Bayram, Husrev T Sencar, and Nasir Memon. 2007. New features to identify computer generated images. In *Image Processing, 2007. ICIP 2007. IEEE International Conference on*, Vol. 4. IEEE, IV–433.

[20] Nhu-Tai Do, In-Seop Na, and Soo-Hyung Kim. 2018. Forensics Face Detection From GANs Using Convolutional Neural Network. In *ISITC’2018*.

[21] Hany Farid. 2009. Image forgery detection. *IEEE Signal processing magazine* 26, 2 (2009), 16–25.

[22] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. 2020. Leveraging Frequency Analysis for Deep Fake Image

Recognition. *arXiv preprint arXiv:2003.08685* (2020).

[23] Jessica Fridrich and Jan Kodovsky. 2012. Rich models for steganalysis of digital images. *IEEE Transactions on Information Forensics and Security* 7, 3 (2012), 868–882.

[24] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.

[25] 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.

[26] Zhiqing Guo, Gaobo Yang, Jiyou Chen, and Xingming Sun. 2020. Fake Face Detection via Adaptive Residuals Extraction Network. *arXiv preprint arXiv:2005.04945* (2020).

[27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Identity mappings in deep residual networks. In *European conference on computer vision*. Springer, 630–645.

[29] Peisong He, Haoliang Li, and Hongxia Wang. 2019. Detection of Fake Images Via The Ensemble of Deep Representations from Multi Color Spaces. In *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2299–2303.

[30] Chih-Chung Hsu, Chia-Yen Lee, and Yi-Xiu Zhuang. 2018. Learning to Detect Fake Face Images in the Wild. In *2018 International Symposium on Computer, Consumer and Control (IS3C)*. IEEE, 388–391.

[31] Nils Hulzebosch, Sarah Ibrahimi, and Marcel Worring. 2020. Detecting CNN-Generated Facial Images in Real-World Scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 642–643.

[32] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. 2017. Globally and locally consistent image completion. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 107.

[33] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-To-Image Translation With Conditional Adversarial Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1125–1134.

[34] Anubhav Jain, Puspita Majumdar, Richa Singh, and Mayank Vatsa. 2020. Detecting GANs and Retouching based Digital Alterations via DAD-HCNN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 672–673.

[35] Anubhav Jain, Richa Singh, and Mayank Vatsa. 2018. On Detecting GANs and Retouching based Synthetic Alterations. In *9th International Conference on Biometrics: Theory, Applications and Systems*.

[36] 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).

[37] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4401–4410.

[38] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2019. Analyzing and Improving the Image Quality of StyleGAN. *arXiv preprint arXiv:1912.04958* (2019).

[39] Jiyeon Kim, Seung-Ah Hong, and Hamin Kim. 2019. A StyleGAN Image Detection Model Based on Convolutional Neural Network. *Journal of Korea Multimedia Society* 22, 12 (2019), 1447–1456.

[40] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. 2017. Learning to Discover Cross-Domain Relations with Generative Adversarial Networks. In *International Conference on Machine Learning*. 1857–1865.

[41] Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[42] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.

[43] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. 2017. Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4681–4690.

[44] Haodong Li, Han Chen, Bin Li, and Shunquan Tan. 2018. Can Forensic Detectors Identify GAN Generated Images?. In *APSIPA Annual Summit and Conference 2018*.

[45] 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).

[46] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang. 2017. Generative face completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3911–3919.

[47] Xiang Lin, Jian-Hua Li, Shi-Lin Wang, Feng Cheng, Xiao-Sa Huang, et al. 2018. Recent Advances in Passive Digital Image Security Forensics: A Brief Review. *Engineering* (2018).

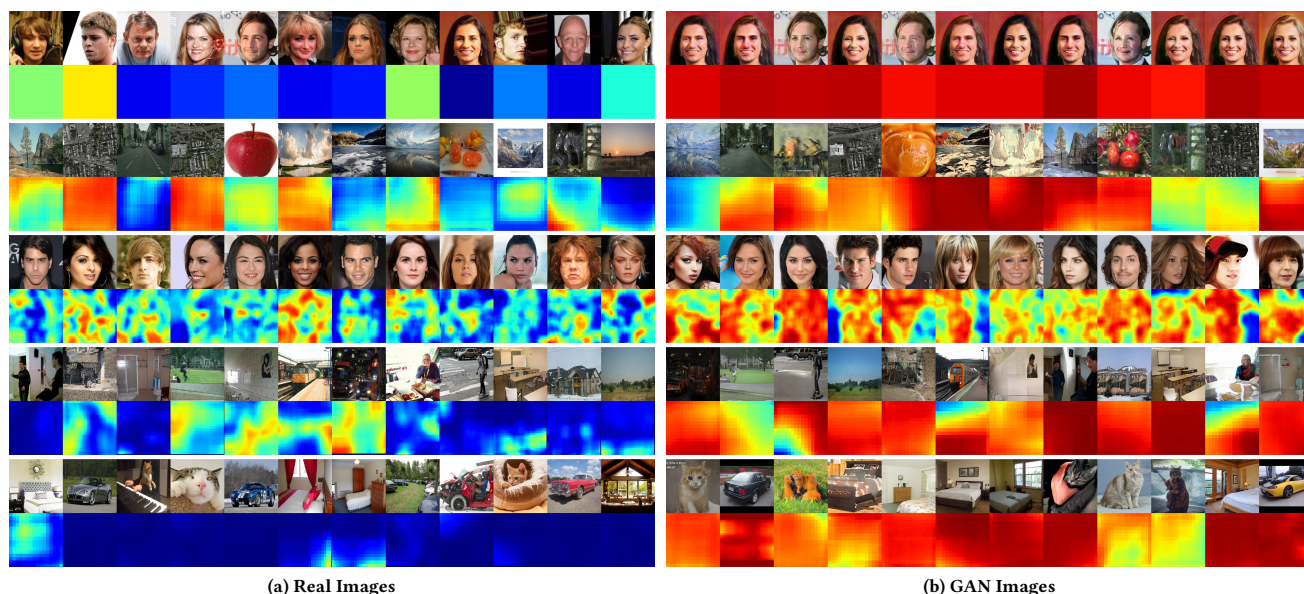


Figure 8: Localization heatmaps of (a) Real images and (b) GAN images from different GAN datasets (top to bottom): ProGAN [36], StarGAN [14], CycleGAN [91], StyleGAN [37], and SPADE/GauGAN [64].

- [48] Ming-Yu Liu and Oncel Tuzel. 2016. Coupled generative adversarial networks. In *Advances in neural information processing systems*. 469–477.
- [49] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*. 3730–3738.
- [50] Zhengzhe Liu, Xiaojuan Qi, and Philip HS Torr. 2020. Global texture enhancement for fake face detection in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8060–8069.
- [51] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, Nov (2008), 2579–2605.
- [52] Brandon Mader, Martin S Banks, and Hany Farid. 2017. Identifying computer-generated portraits: The importance of training and incentives. *Perception* 46, 9 (2017), 1062–1076.
- [53] Babak Mahdian and Stanislav Saic. 2010. A bibliography on blind methods for identifying image forgery. *Signal Processing: Image Communication* 25, 6 (2010), 389–399.
- [54] Hadi Mansourifar and Weidong Shi. 2020. One-Shot GAN Generated Fake Face Detection. *arXiv preprint arXiv:2003.12244* (2020).
- [55] Francesco Marra, Diego Gagnaniello, Davide Cozzolino, and Luisa Verdoliva. 2018. Detection of GAN-generated fake images over social networks. In *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*. IEEE, 384–389.
- [56] Francesco Marra, Diego Gagnaniello, Luisa Verdoliva, and Giovanni Poggi. 2018. Do GANs leave artificial fingerprints? *arXiv preprint arXiv:1812.11842* (2018).
- [57] Francesco Marra, Cristiano Saltori, Giulia Boato, and Luisa Verdoliva. 2019. Incremental learning for the detection and classification of GAN-generated images. *arXiv preprint arXiv:1910.01568* (2019).
- [58] Scott McCloskey and Michael Albright. 2018. Detecting GAN-generated Imagery using Color Cues. *arXiv preprint arXiv:1812.08247* (2018).
- [59] Mehdi Mirza and Simon Osindero. 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784* (2014).
- [60] 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.
- [61] Lakshmanan Nataraj, Tajuddin Manhar Mohammed, BS Manjunath, Shivkumar Chandrasekaran, Arjuna Flenner, Jawadul H Bappy, and Amit K Roy-Chowdhury. 2019. Detecting GAN generated fake images using co-occurrence matrices. In *Media Watermarking, Security, and Forensics*. IS&T International Symposium on Electronic Imaging.
- [62] J Neves, Ruben Tolosana, Ruben Vera-Rodriguez, Vasco Lopes, Hugo Proença, and Julian Fierrez. 2019. GANprintR: Improved Fakes and Evaluation of the State-of-the-Art in Face Manipulation Detection. *arXiv preprint arXiv:1911.05351* (2019).
- [63] Joao C Neves, Ruben Tolosana, Ruben Vera-Rodriguez, Vasco Lopes, and Hugo Proença. 2019. Real or Fake? Spoofing State-Of-The-Art Face Synthesis Detection Systems. *arXiv preprint arXiv:1911.05351* (2019).
- [64] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. 2019. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2337–2346.
- [65] Guim Perarnau, Joost van de Weijer, Bogdan Raducanu, and Jose M Álvarez. 2016. Invertible conditional gans for image editing. *arXiv preprint arXiv:1611.06355* (2016).
- [66] Tomáš Pevný, Patrick Bas, and Jessica Fridrich. 2010. Steganalysis by subtractive pixel adjacency matrix. *Information Forensics and Security, IEEE Transactions on* 5, 2 (2010), 215–224.
- [67] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434* (2015).
- [68] Nicolas Rahmouni, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. 2017. Distinguishing computer graphics from natural images using convolution neural networks. In *Information Forensics and Security (WIFS), 2017 IEEE Workshop on*. IEEE, 1–6.
- [69] Yuan Rao and Jiangqun Ni. 2016. A deep learning approach to detection of splicing and copy-move forgeries in images. In *Information Forensics and Security (WIFS), 2016 IEEE International Workshop on*. IEEE, 1–6.
- [70] 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.
- [71] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [72] Kenneth Sullivan, Upamanyu Madhow, Shivkumar Chandrasekaran, and BS Manjunath. 2006. Steganalysis for Markov cover data with applications to images. *IEEE Transactions on Information Forensics and Security* 1, 2 (2006), 275–287.
- [73] Kenneth Sullivan, Upamanyu Madhow, Shivkumar Chandrasekaran, and Bangalore S Manjunath. 2005. Steganalysis of spread spectrum data hiding exploiting cover memory. In *Security, Steganography, and Watermarking of Multimedia Contents VII*, Vol. 5681. International Society for Optics and Photonics, 38–47.
- [74] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. 2017. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- [75] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.

- [76] Yaniv Taigman, Adam Polyak, and Lior Wolf. 2016. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200* (2016).
- [77] 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.
- [78] Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-Garcia. 2020. Deepfakes and beyond: A survey of face manipulation and fake detection. *arXiv preprint arXiv:2001.00179* (2020).
- [79] Rafael Valle, UC CNMAT, Wilson Cai, and Anish Doshi. 2018. TequilaGAN: How to easily identify GAN samples. *arXiv preprint arXiv:1807.04919* (2018).
- [80] Savita Walia and Krishan Kumar. 2018. Digital image forgery detection: a systematic scrutiny. *Australian Journal of Forensic Sciences* (2018), 1–39.
- [81] Run Wang, Lei Ma, Felix Juefei-Xu, Xiaofei Xie, Jian Wang, and Yang Liu. 2019. FakeSpotter: A Simple Baseline for Spotting AI-Synthesized Fake Faces. *arXiv preprint arXiv:1909.06122* (2019).
- [82] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. 2019. CNN-generated images are surprisingly easy to spot... for now. *arXiv preprint arXiv:1912.11035* (2019).
- [83] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8798–8807.
- [84] Ruoyu Wu, Xiaolong Li, and Bin Yang. 2011. Identifying computer generated graphics via histogram features. In *Image Processing (ICIP), 2011 18th IEEE International Conference on*. IEEE, 1933–1936.
- [85] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. 2017. DualGAN: Unsupervised Dual Learning for Image-to-Image Translation. In *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2868–2876.
- [86] Ning Yu, Larry Davis, and Mario Fritz. 2018. Attributing fake images to GANs: Analyzing fingerprints in generated images. *arXiv preprint arXiv:1811.08180* (2018).
- [87] Kejun Zhang, Yu Liang, Jianyi Zhang, Zhiqiang Wang, and Xinxin Li. 2019. No One Can Escape: A General Approach to Detect Tampered and Generated Image. *IEEE Access* 7 (2019), 129494–129503.
- [88] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. 2019. Detecting and Simulating Artifacts in GAN Fake Images. *arXiv preprint arXiv:1907.06515* (2019).
- [89] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. 2017. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 633–641.
- [90] Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. 2018. Learning Rich Features for Image Manipulation Detection. *arXiv preprint arXiv:1805.04953* (2018).
- [91] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In *IEEE International Conference on Computer Vision*.
- [92] Yi-Xiu Zhuang and Chih-Chung Hsu. 2019. Detecting Generated Image Based on a Coupled Network with Two-Step Pairwise Learning. In *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 3212–3216.