

High-Resolution Deep Convolutional Generative Adversarial Networks

J. D. Curtó^{*,1,2,3,4}, I. C. Zarza^{*,1,2,3,4}, F. De La Torre²,
I. King¹, and M. R. Lyu¹

¹Dept. of Computer Science and Engineering, The Chinese University of Hong Kong

²The Robotics Institute, Carnegie Mellon

³Laboratory of Computer Vision, Eidgenössische Technische Hochschule Zürich

⁴Dept. of Electronic Engineering, City University of Hong Kong

{curto,zarza,king,lyu}@cse.cuhk.edu.hk, ftorre@cs.cmu.edu

*Both authors contributed equally

Abstract. Generative Adversarial Networks (GANs) [1] convergence in a high-resolution setting with a computational constrain of GPU memory capacity (from 12GB to 24 GB) has been beset with difficulty due to the known lack of convergence rate stability. In order to boost network convergence of DCGAN (Deep Convolutional Generative Adversarial Networks) [2] and achieve good-looking high-resolution results we propose a new layered network structure, HDCGAN, that incorporates current state-of-the-art techniques for this effect. A novel dataset containing human faces from different ethnical groups in a wide variety of illumination conditions and image resolutions is introduced, Curtó & Zarza¹. Curtó is enhanced with HDCGAN synthetic images, thus being the first GAN augmented face dataset. We conduct extensive experiments on CelebA [3] (MS-SSIM 0.1978 and Distance of Fréchet 8.77) and Curtó.

1 Introduction

Developing a Generative Adversarial Network (GAN) structure [1] able to produce good quality high-resolution samples from images has important applications including image inpainting, 3D data, localization and semi-supervised learning.

In this paper, we focus on the task of face generation, as it gives GANs a huge space of learning attributes. In this context, we introduce the Dataset of Curtó & Zarza, a well-balanced collection of images containing 14,248 human faces from different ethnical groups and rich in a wide range of learnable attributes, such as gender and age diversity, hair-style and pose variation or presence of smile, glasses, hats and fashion items. We also ensure the presence of changes in illumination and image resolution. We propose to use Curtó as de facto approach to empirically test the distribution learned by a GAN, as it offers a challenging problem to solve, while keeping the number of samples, and therefore training time, bounded. A set of random samples can be seen in Figure 1.

¹ Curtó is available at <https://www.github.com/curto2/c/>

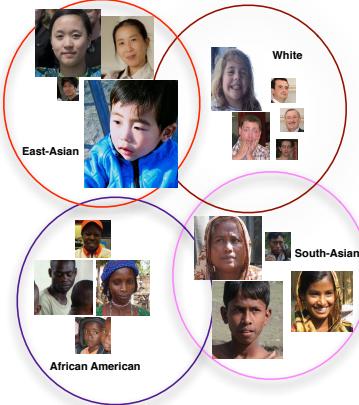


Figure 1. Image Samples of Curtó. A set of random samples for each ethnicity category.

Despite improvements in GANs training stability [4] and specific-task design during the last years, it is still challenging to train GANs to generate high-resolution images due to the disjunction in the high dimensional pixel space between supports of the real image and implied model distributions [5], [6].

Our goal is to be able to generate indistinguishable sample instances using face data to push the boundaries of GAN image generation that scale well to high-resolution images (i.e. 512×512) and where context information is maintained.

In this sense, Deep Learning has a tremendous appetite for data. The question that arises instantly is, what if we were able to generate additional realistic data to aid the learning process using the same techniques that are later used to train the system. The first step would then be to have an image generation tool able to sample from a very precise distribution (e.g. faces from celebrities) which instances resemble or highly correlate with real sample images of the underlying true distribution. Once achieved, what is desirable and comes next is that these generated image points not only fit well into the original distribution set of images but also add additional useful information such as redundancy, different poses or even generate highly-probable scenarios that would be possible to see in the original dataset but are actually not present.

To achieve the former goal this work contributes in the following:

- Network structure that achieves compelling results and scales well to the high-resolution setting where to the best of our knowledge other variant network architectures are unable to continue learning or fall into mode collapse.

- New dataset targeted for GAN training, Curtó, that introduces a wide space of learning attributes. It aims to provide a well-posed difficult task while keeping training time and resources tightly bounded to spearhead research in the area.

2 Prior Work

Generative image generation is a key problem in Computer Vision. Remarkable advances have been made with the renaissance of Deep Learning. Variational Autoencoders (VAE) [7] formulates the problem with a probabilistic graphical model approach, where the lower bound of data likelihood is maximized. Autoregressive models (i.e. PixelRNN [8]), based on modeling the conditional distribution of the pixel space, have also presented relative success generating synthetic images. Lately, Generative Adversarial Networks (GANs) [1], [2], [9], [10], [11], [12] have shown strong performance in image generation. However, training instability makes it very hard to scale to high-resolution (256×256 or 512×512) samples. Some current works on the topic pinpoint this specific problem [13], where conditional image generation is also tackled while other recent techniques [4], [14], [15], [16], [17] try to stabilize the training process.

3 Dataset of Curtó & Zarza

Curtó contains 14,248 face images balanced in terms of ethnicity: african american, east-asian, south-asian and white. Mirror images are included to enhance pose variation and there is roughly 25% per image class. Labels consist on JSON files with thorough attribute information: gender, age, ethnicity, hair color, hair style, eyes color, facial hair, glasses, visible forehead, hair covered and smile. There is also an extra set with 3,384 cropped labeled face images, ethnicity white, no mirror samples included. We crawled Flickr to download face images from several countries that contain different hair-style variations and style attributes. These images were then processed to extract 49 facial landmark points using [18]. We ensure using Amazon Mechanical Turk that the detected faces are correct in terms of ethnicity and face detection. Cropped faces are then extracted to generate multiple resolution sources. Mirror augmentation is performed to further enhance pose variation.

Curtó introduces a difficult learning paradigm, where different ethnical groups are present, with very varied fashion and hair styles. The fact that the photos are taken using non-professional cameras in a non-controlled environment, gives us multiple face poses, illumination conditions and camera quality.

4 Approach

Generative Adversarial Networks (GANs) proposed by [1] are based on two dueling networks; Generator G and Discriminator D . In essence, the learning process

consists on a two-player game where D tries to distinguish between the prediction of G and the ground truth, while at the same time G tries to fool D by producing fake instance samples as closer to the real ones as possible.

The min-max game entails the following objective function.

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))] . \quad (1)$$

where x is a ground truth image sampled from the true distribution p_{data} , and z is a noise vector sampled from p_z (i.e. uniform or normal distribution). G and D are parametric functions where $G : p_z \rightarrow p_{data}$ maps samples from noise distribution p_z to data distribution p_{data} .

As an extension to this framework, DCGAN [2] proposes an architectural topology based on Convolutional Neural Networks (CNNs) to stabilize training and re-use state-of-the-art network topologies from image classification tasks. This direction has recently received lots of attention due to its compelling results in several supervised and unsupervised learning problems. We build on this to propose a novel DCGAN architecture to address the problem of high-resolution image generation. We name this approach HDCGAN.

4.1 HDCGAN

Despite the undoubtable success, GANs are still arduous to train, particularly when we use big images (e.g. 512×512). It is very common to see D beating G in the learning process, or the reverse, ending in unrecognizable imagery, also known as mode collapse. Just when stable learning is achieved, the GAN structure is able to succeed in getting better and better results with time.

This issue is what drives us to carefully derive a simple yet powerful structure that leverages common problems and gets a stable and steady training mechanism.

Self-normalizing Neural Networks (SNNs) were introduced in [19]. We consider a neural network with activation function f , connected to the next layer by a weight matrix \mathbf{W} , and whose inputs are the activations from the preceding layer x , i.e. $y = f(\mathbf{W}x)$.

We can define a mapping g that maps mean and variance from one layer to mean and variance of the following layer

$$\begin{pmatrix} \mu \\ \nu \end{pmatrix} \longmapsto \begin{pmatrix} \tilde{\mu} \\ \tilde{\nu} \end{pmatrix} : \begin{pmatrix} \tilde{\mu} \\ \tilde{\nu} \end{pmatrix} = g \begin{pmatrix} \mu \\ \nu \end{pmatrix} . \quad (2)$$

Common normalization tactics such as batch normalization ensure a mapping g that keeps (μ, ν) and $(\tilde{\mu}, \tilde{\nu})$ close to a desired value, normally $(0, 1)$.

SNNs go beyond this assumption and require the existence of a mapping $g : \Omega \rightarrow \Omega$ that for each activation y maps mean and variance from one layer to the next layer and at the same time have a stable and attracting fixed point depending on (ω, τ) in Ω . Moreover, the mean and variance remain in the domain Ω and when iteratively applying the mapping g , each point within Ω converges to this fixed point. Therefore, SNNs keep activations normalized when propagating them through the layers of the network.

Scaled Exponential Linear Units (SELU) [19] is introduced as the choice of activation function in Feed-forward Neural Networks (FNNs) to construct a mapping g with properties that lead to SNNs.

$$selu(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha \exp^x - \alpha & \text{if } x \leq 0. \end{cases} \quad (3)$$

Empirical observation leads us to say that the use of SELU activation function greatly improves the convergence speed on the DCGAN structure, however, after some iterations mode collapse and gradient explosion completely destroy training when using high-resolution images. We conclude that although SELU gives theoretical guarantees as the optimal activation function in FNNs, numerical errors in the GPU computation degrade its performance in the overall min-max game of DCGAN. To alleviate this problem, we propose to use SELU and BatchNorm [20] together. The motivation is that when numerical errors move $(\tilde{\mu}, \tilde{\nu})$ away from the attracting point that depends on $(\omega, \tau) \in \Omega$, BatchNorm will ensure it is close to a desired value and therefore maintain the convergence rate. Experiments show that this technique stabilizes training and allows us to use fewer GPU resources, having steady diminishing errors in G and D . It also accelerates convergence speed by a great factor, as can be seen after just some few epochs of CelebA training in Figure 5.

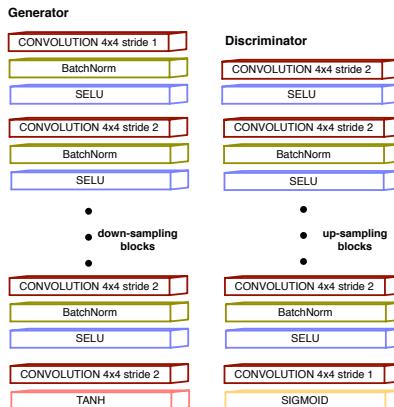


Figure 2. HDCGAN Architecture. Generator and Discriminator.

HDCGAN Architecture is described in Figure 2. It differs from traditional DCGAN in the use of SELU + BatchNorm (BS) layers instead of ReLUs.

We observe that when having difficulty in training DCGAN, it is always better to use a fixed learning rate and instead increase the batch size. This is because having more diversity in training, gives a steady diminishing loss and better generalization results. To aid the learning process, additive noise is added to both the inputs of D and G . We see that this helps overcome mode saturation and collapse.

We empirically show that the use of BS induces SNNs properties in the GAN structure, and thus makes the learning procedure highly robust, even in the stark presence of noise and perturbations. This behavior can be observed when the zero-sum game problem stabilizes and errors in D and G jointly diminish, Figure 6. Comparison to traditional DCGAN, Wasserstein GAN [21] and WGAN-GP [22] is not possible, as to date, all other former methods, such as [23], cannot generate recognizable results in image size 512×512 , 24GB GPU memory setting.

Thus, HDCGAN pushes up state-of-the-art results beating all former DCGAN-based architectures and shows that, under the right circumstances, BS can solve the min-max game efficiently.

5 Empirical Analysis

We build on DCGAN and extend the framework to train with high-resolution images using Pytorch. Our experiments are conducted using a fixed learning rate of 0.0002 and ADAM solver [24] with batch size 32 and 512×512 training images with the number of filters of G and D equal to 64.

In order to test generalization capability, we train HDCGAN in the newly introduced Curtó and CelebA.

Technical Specifications: $2 \times$ NVIDIA Titan X, Intel Core i7-5820k@3.30GHz.

5.1 Curtó

The results after 150 epochs are shown in Figure 3. We can see that HDCGAN captures the underlying image features that represent faces and not only memorizes training examples. We retrieve nearest neighbors to the generated images in Figure 4 to illustrate this effect.



Figure 3. HDCGAN Example Results, Dataset of Curtó & Zarza. 150 epochs of training. Image size 512×512 .



Figure 4. Nearest Neighbors. Generated images in the first row and retrieving their five nearest neighbors in the training images (rows 2-6).

5.2 CelebA

CelebA is a large-scale dataset with 202,599 celebrity faces. It mainly contains frontal faces and is particularly biased towards white ethnical groups. The fact that it presents very controlled illumination settings and good photo resolution, makes it a considerably easier problem than Curtó. The results after just 12 epochs of training are shown in Figure 5.



Figure 5. HDCGAN Example Results, CelebA Dataset. 12 epochs of training. Image size 512×512.

In Figure 6 we can observe that BS stabilizes the zero-sum game problem. To show the validity of our method, we enclose Figure 7 and Figure 10, presenting a large number of samples for epochs 19 and 39 of the learning process. We also attach zoomed-in examples to appreciate the quality and size of the generated samples, Figure 8 and Figure 11. Failure cases can be observed in Figure 9 and Figure 12.

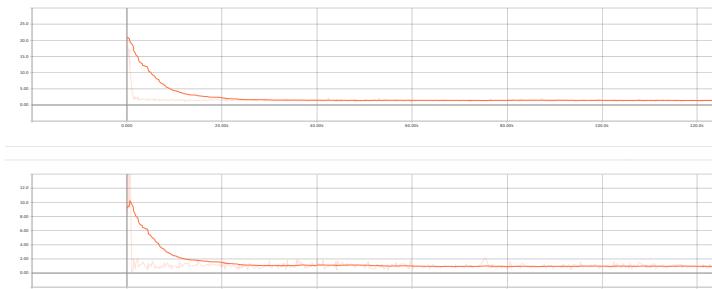


Figure 6. HDCGAN on CelebA. Error in Discriminator (top) and Error in Generator (bottom). 19 epochs of training.

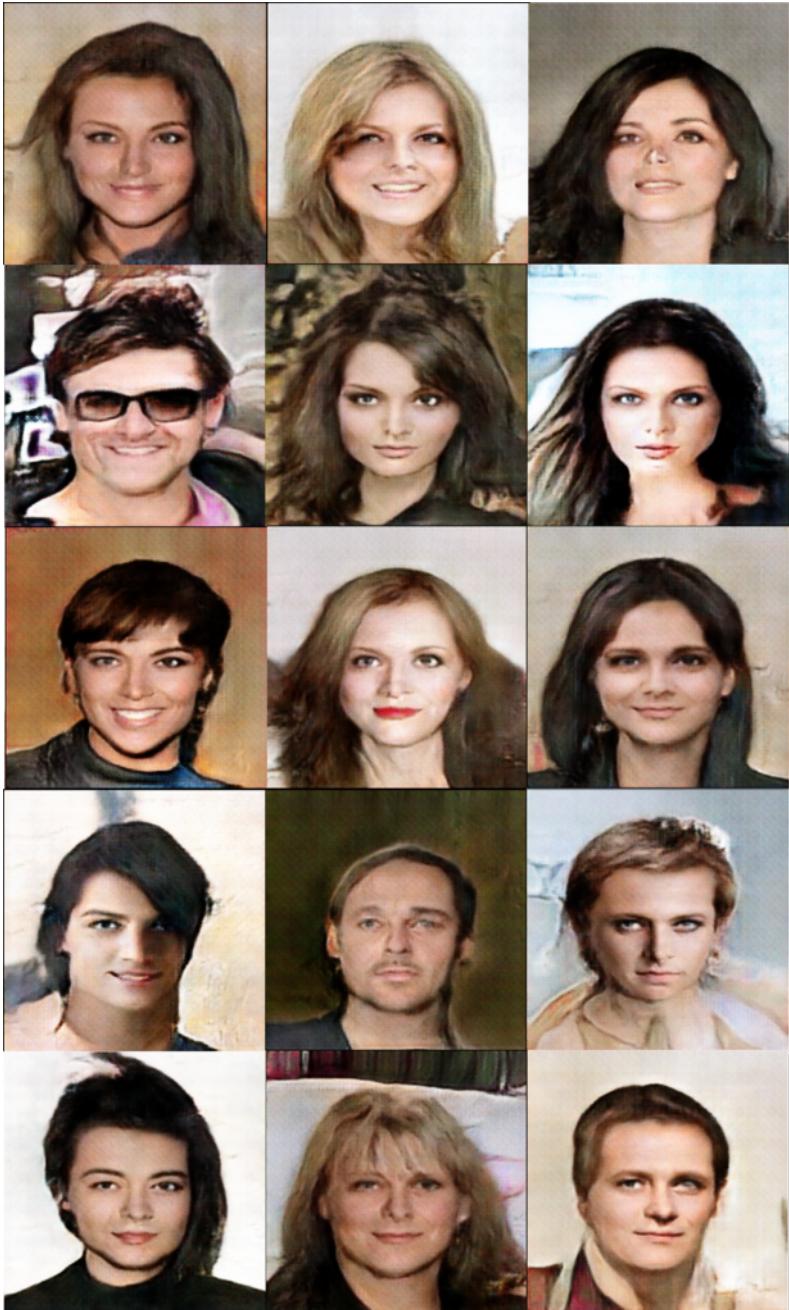


Figure 7. HDCGAN Example Results, CelebA Dataset. 19 epochs of training.
Image size 512×512 .

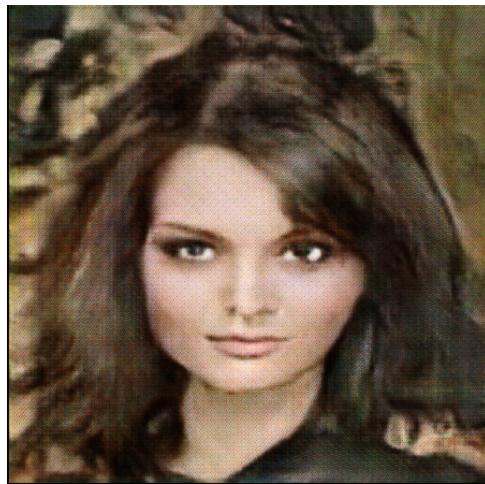


Figure 8. HDCGAN Example Result, CelebA Dataset. 19 epochs of training.
Image size 512×512. 35% of full-scale image.



Figure 9. HDCGAN Example Results, CelebA Dataset. 19 epochs of training.
Image size 512×512. Failure cases.



Figure 10. HDCGAN Example Results, CelebA Dataset. 39 epochs of training.
Image size 512×512.

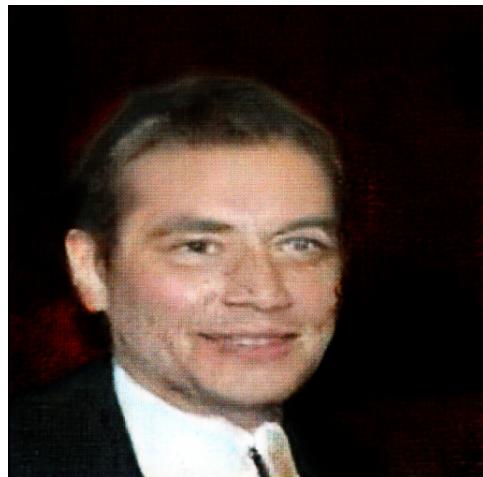


Figure 11. HDCGAN Example Result, CelebA Dataset. 39 epochs of training.
Image size 512×512. 35% of full-scale image.



Figure 12. HDCGAN Example Results, CelebA Dataset. 39 epochs of training.
Image size 512×512. Failure cases.

Besides, to illustrate how fundamental our approach is, we enlarge Curtó with 4,239 unlabeled synthetic images generated by HDCGAN on CelebA, a random set can be seen in Figure 13.

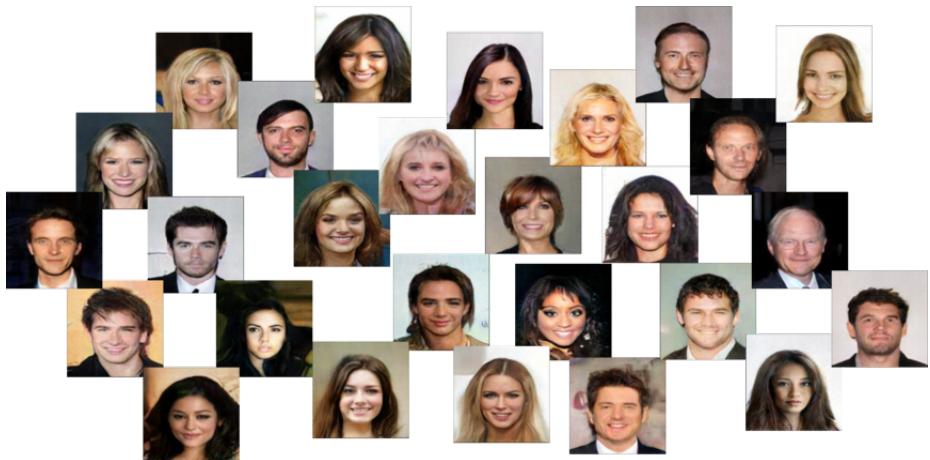


Figure 13. HDCGAN Synthetic Images. A set of random samples.

6 Assessing the Discriminability and Quality of Generated Samples

We build on previous image similarity metrics to qualitatively evaluate generated samples of Generative Models. The most effective of these is multi-scale structural similarity (MS-SSIM) [9]. We make comparison at resized image size 128×128 on CelebA. MS-SSIM results are averaged from 10,000 pairs of generated images. Table 1 shows HDCGAN significantly improves state-of-the-art results.

Table 1. Multi-scale structural similarity (MS-SSIM) results on CelebA at resized image size 128×128 . Lower is better.

	MS-SSIM
Gulrajani et al. (2017) [22]	0.2854
Karras et al. (2018) [17]	0.2838
HDCGAN	0.1978

We monitor MS-SSIM scores across several training epochs averaging from 10,000 pairs of generated images to see the temporal performance, Figure 14. HDCGAN improves the quality of the samples while increases the diversity of the generated distribution.

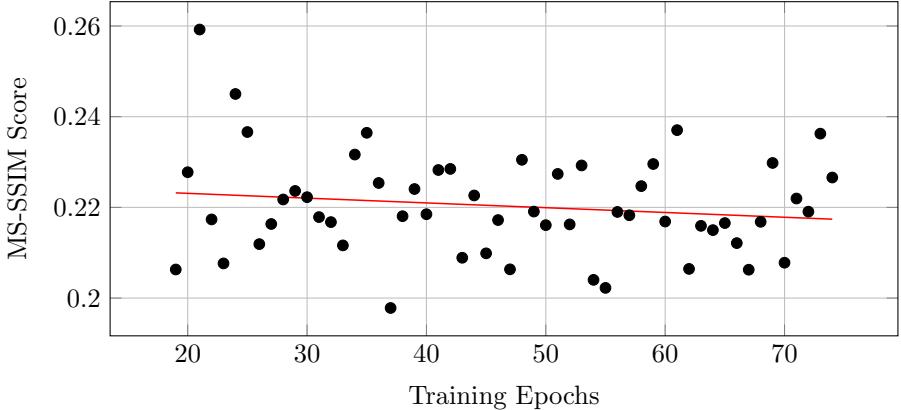


Figure 14. MS-SSIM Scores on CelebA across several training epochs. Results are averaged from 10,000 pairs of generated images from training epoch 19 to 74. Comparison is made at resized image size 128×128 . Affine interpolation is shown in red.

In [25] they propose to evaluate GANs using the distance of Fréchet, which assesses the similarity between two distributions by the difference of two Gaussians. We make comparison at resized image size 64×64 on CelebA. Results are computed from 10,000 512×512 generated samples from epochs 36 to 52, resized at image size 64×64 yielding a value of 8.77.

7 Discussion

In this paper, we propose High-Resolution Deep Convolutional Generative Adversarial Networks (HDCGAN) by stacking SELU + BatchNorm (BS) layers. The proposed method generates high-resolution images (e.g. 512×512) in circumstances where all other former methods fail. It exhibits a steady and smooth training mechanism. HDCGAN is the current state-of-the-art in synthetic image generation on CelebA (MS-SSIM 0.1978 and Distance of Fréchet 8.77).

Further, we present a face dataset containing well-balanced ethnical groups for GAN training, Curtó & Zarza, that poses a very difficult challenge and is rich on learning attributes to sample from. Moreover, we enhance Curtó with 4,239 unlabeled synthetic images generated by HDCGAN, being therefore the first GAN augmented face dataset.

References

- [1] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial neworks. In: NIPS. (2014) [1](#), [3](#)
- [2] Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial network. In: ICLR. (2016) [1](#), [3](#), [4](#)
- [3] Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: ICCV. (2015) [1](#)
- [4] Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Redford, A., Chen, X.: Improved techniques for training gans. In: NIPS. (2016) [2](#), [3](#)
- [5] Arjovsky, M., Bottou, L.: Towards principled methods for training generative adversarial networks. In: ICLR. (2017) [2](#)
- [6] Sønderby, C., Caballero, J., Theis, L., Shi, W., Hussar, F.: Amortised map inference for image super-resolution. In: ICLR. (2017) [2](#)
- [7] Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: ICLR. (2014) [3](#)
- [8] van den Oord, A., Kalchbrenner, N., Kavukcuoglu, K.: Pixel recurrent neural networks. In: ICML. (2016) [3](#)
- [9] Odena, A., Olah, C., Shlens, J.: Conditional image synthesis with auxiliary classifier gans. In: ICML. (2017) [3](#), [13](#)
- [10] Antoniou, A., Storkey, A., Edwards, H.: Data augmentation generative adversarial networks. In: ICLR. (2018) [3](#)
- [11] Wang, X., Gupta, A.: Generative image modeling using style and structure adversarial networks. In: ECCV. (2016) [3](#)
- [12] Zhu, J., Krähenbühl, P., Shechtman, E., Efros, A.A.: Generative visual manipulation on the natural image manifold. In: ECCV. (2016) [3](#)
- [13] Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., Metaxas, D.: Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In: CVPR. (2017) [3](#)
- [14] Chen, Q., Koltun, V.: Photographic image synthesis with cascaded refinement networks. In: ICCV. (2017) [3](#)
- [15] Dosovitskiy, A., Brox, T.: Generating images with perceptual similarity metrics based on deep networks. In: NIPS. (2016) [3](#)
- [16] Zhao, J., Mathieu, M., LeCun, Y.: Energy-based generative adversarial network. In: ICLR. (2017) [3](#)
- [17] Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. In: ICLR. (2018) [3](#), [13](#)
- [18] Xiong, X., De La Torre, F.: Supervised descent method and its application to face alignment. In: CVPR. (2013) [3](#)
- [19] Klarbauer, G., Unterthiner, T., Mayr, A.: Self-normalizing neural networks. In: NIPS. (2017) [4](#), [5](#)
- [20] Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML. (2015) [5](#)

- [21] Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein gan. In: ICML. (2017) 6
- [22] Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.: Improved training of wasserstein gans. In: NIPS. (2017) 6, 13
- [23] Denton, E., Chintala, S., Szlam, A., Fergus, R.: Deep generative image models using a laplacian pyramid of adversarial networks. In: NIPS. (2015) 6
- [24] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR. (2015) 6
- [25] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: NIPS. (2017) 14