# Semi-supervised Adversarial Learning to Generate Photorealistic Face Images of New Identities from 3D Morphable Model

Baris Gecer[1], Binod Bhattarai[1], Josef Kittler[2], and Tae-Kyun Kim[1]

[1]Department of Electrical and Electronic Engineering, Imperial College London, UK
[2]Centre for Vision, Speech and Signal Processing, University of Surrey, UK
{b.gecer,b.bhattarai,tk.kim}@imperial.ac.uk,
j.kittler@surrey.ac.uk

**Abstract.** We propose a novel end-to-end semi-supervised adversarial framework to generate photorealistic face images of new identities with wide ranges of expressions, poses, and illuminations conditioned by a 3D morphable model. Previous adversarial style-transfer methods either supervise their networks with large volume of paired data or use unpaired data with a highly under-constrained two-way generative framework in an unsupervised fashion. We introduce pairwise adversarial supervision to constrain two-way domain adaptation by a small number of paired real and synthetic images for training along with the large volume of unpaired data. Extensive qualitative and quantitative experiments are performed to validate our idea. Generated face images of new identities contain pose, lighting and expression diversity and qualitative results show that they are highly constraint by the synthetic input image while adding photorealism and retaining identity information. We combine face images generated by the proposed method with the real data set to train face recognition algorithms. We evaluated the model on two challenging data sets: LFW and IJB-A. We observe that the generated images from our framework consistently improves over the performance of deep face recognition network trained with Oxford VGG Face dataset and achieves comparable results to the state-of-the-art.

## 1 Introduction

Deep learning has shown an great improvement in performance of several computer vision tasks [37,19,15,11,12,58] including face recognition [33,41,54,30,53] in the recent years. This was mainly thanks to the availability of large-scale datasets. Yet the performance is often limited by the volume and the variations of training examples. Larger and wider datasets usually improve the generalization and overall performance of the model [41,1].

The process of collecting and annotating training examples for every specific computer vision task is laborious and non-trivial. To overcome this challenge, additional synthetic training examples along with limited real training examples can be utilised to train the model. Some of the recent works such as 3D face reconstruction [38], gaze estimation [61,52], human pose, shape and motion estimation [49] *etc.* use additional
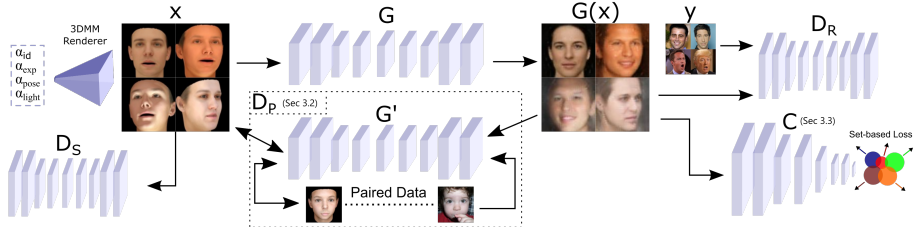
Fig. 1: Our approach aims to synthesize photorealistic images conditioned by a given synthetic image by 3DMM. It regularizes cycle consistency [63] by introducing an additional adversarial game between the two generator networks in an unsupervised fashion. Thus the under-constraint cycle loss is supervised to have correct matching between the two domains by the help of a limited number of paired data. We also encourage the generator to preserve face identity by a set-based supervision through a pretrained classification network.

synthetic images generated from 3D models to train deep networks. One can generate synthetic face images using a 3D morphable model (3DMM) [3] by manipulating identity, expression, illumination, and pose parameters. However, the resulting images are not photorealistic enough to be suitable for in-the-wild face recognition tasks. It is beacause the information of real face scans is compressed by the 3DMM and the graphical engine that models illumination and surface is not perfectly accurate. Thus, the main challenge of using synthetic data obtained from 3DMM model is the discrepancy in nature and quality of synthetic and real images which pose the problem of domain adaptation [34]. Recently, adversarial training methods [42,44,10] become popular to mitigate such challenges.

Generative Adversarial Network (GAN), introduced by Goodfellow *et al*. [17], and its variants [35,24,2,13] are quite successful in generating realistic images. However, in practice, GANs are likely to stuck in mode collapse for large scale image generation. They are also unable to produce images that are 3D coherent and globally consistent [17]. To overcome these drawbacks, we propose a semi-supervised adversarial learning framework to synthesize photorealistic face images of new identities with numerous data variation supplied by a 3DMM. We address these shortcomings by exciting a generator network with synthetic images sampled from 3DMM and transforming them into photorealistic domain using adversarial training as a bridge. Unlike most of the existing works that excite their generators with a noise vector [35,2], we feed our generator network by synthetic face images. Such a strong constraint naturally helps in avoiding the mode collapse problem, one of the main challenges faced by the current GAN methods. Fig. 1 shows the general overview of the proposed method. We discuss the proposed method in more details in Sec. 3.

In this paper, we address the challenge of generating photorealistic face images from 3DMM rendered faces of different identities with arbitrary poses, expressions, and illuminations. We formulate this problem as a domain adaptation problem *i.e.* aligning the 3DMM rendered face domain into realistic face domain. One of the previous works

closest to ours [22] address style transfer problem between a pair of domains with classical conditional GAN. The major bottleneck of this method is, it requires a large number of paired examples from both domains which are hard to collect. CycleGAN [63], another recent method and closest to our work, proposes a two-way GAN framework for unsupervised image-to-image translation. However, the cycle consistency loss proposed in their method is satisfied as long as the transitivity of the two mapping networks is maintained. Thus, the resulting mapping is not guaranteed to produce the intended transformation. To overcome the drawbacks of these methods [22,63], we propose to use a small amount of paired data to train an inverse mapping network as a matching aware discriminator. In the proposed method, the inverse mapping network plays the role of both the generator and the discriminator. To the best of our knowledge, this is the first attempt for adversarial semi-supervised style translation for an application with such limited paired data.

Adding realism to the synthetic face images and preserving their identity information is a challenging problem. Although synthetic input images, 3DMM rendered faces, contain distinct face identities, the distinction between them vanishes as a result of the virtue of non-linear transformations while the discriminator encourages realism. To tackle such problem, prior works either employ a separate pre-trained network [57] or embed Identity labels (id) [46] into the discriminator. Unlike existing works, which are focused on generating new images of existing identities, we are interested in generating multiple images of new identities itself. Therefore, such techniques are not directly applicable to our problem. To address this challenge, we propose to use set-based center [50] and pushing loss functions [16] on top of a pre-trained face embedding network. This will keep track of the changing average of embeddings of generated images belonging to same identity (i.e. centroids). In this way identity preservation becomes adaptive to changing feature space during the training of the generator network unlike softmax layer that converges very quickly at the beginning of the training before meaningful images are generated.

Our contributions can be summarized as follows:

– We propose a novel end-to-end adversarial training framework to generate photorealistic face images of new identities constrained by synthetic 3DMM images with identity, pose, illumination and expression diversity. The resulting synthetic face images are visually plausible and can be used to boost face recognition as additional training data or any other graphical purposes.
– We propose a novel semi-supervised adversarial style transfer approach that trains an inverse mapping network as a discriminator with paired synthetic-real images.
– We employ a novel set-based loss function to preserve consistency among unknown identities during GAN training.

## 2   Related Works

In this Section we discuss the prior art that is closely related to the proposed method.

*Domain Adaptation.* As stated in the introduction, our problem of generating photorealistic face images from 3DMM rendered faces can be seen as a domain adaptation

problem. A straightforward adaptation approach is to align the distributions at the feature level by simply adding a loss to measure the mismatch either through second-order moments [45] or with adversarial losses [47,48,14].

Recently, pixel level domain adaptation becomes popular due to practical breakthroughs on Kullback-Leibler divergence [18,17,35], namely GANs which optimize a generative and discriminative network through a mini-max game. It has been applied to a wide range of applications including fashion clothing [27], person specific avatar creation [51], text-to-image synthesis [59], face frontalization [57], and retinal image synthesis [10].

Pixel domain adaptation can be done in a supervised manner simply by conditioning the discriminator network [22] or directly the output of the generator [8] with the expected output when there is enough paired data from both domains. Please note collecting a large number of paired training examples is expensive, and often requires expert knowledge. [36] proposes a text-to-image synthesis GAN with a matching aware discriminator. They optimize their discriminator for image-text matching beside requiring realism with an additional mismatched text-image pair.

For the cases where paired data is not available, many approaches take an unsupervised way such as pixel-level consistency between input and output of the generator network [5,42], an encoder architecture that is shared by both domains[6] and adaptive instance normalization [21]. An interesting approach is to have two way translation between domains with two distinct generator and discriminator networks. They constrain the two mappings to be inverses of each other with either ResNet [63] or encoder-decoder network [29] as the generator.

*Synthetic Training Data Generation.* The usage of synthetic data as additional training data is shown to be helpful even if they are graphically rendered images in many applications such as 3D face reconstruction [38], gaze estimation [61,52], human pose, shape and motion estimation [49]. Despite the availability of almost infinite number of synthetic images, those approaches are limited due to the domain difference from that of in-the-wild images.

Many existing works utilized adversarial domain adaptation to translate images into photorealistic domain such that they are more useful as a training data. [62] generates many unlabeled samples to improve person re-identification in a semi-supervised fashion. RenderGAN [44] proposes a sophisticated approach to refine graphically rendered synthetic images of tagged bees to be used as training data for bee tag decoding application. WaterGAN [28] synthesizes realistic underwater images by modeling camera parameters and environment effects explicitly to be used as training data for color correction task. Some studies deform existing images by a 3D model to augment diverse set of dataset [32] without adversarial learning.

One of the recent works, simGAN [42], generates realistic synthetic data to improve eye gaze and hand pose estimation. It optimizes pixel level correspondence between input and output of the generator network to preserve content of the synthetic image. This is in fact a limited solution since the pixel-consistency loss encourages the generated images to be similar to synthetic input images and it partially contradicts adversarial realism loss. Instead, we employ an inverse translation network similar to cycleGAN [63] with an additional pair-wise supervision to preserve the initial condition without hurt-

ing realism. This network also behaves as a discriminator to a straight mapping network with a real paired data to avoid possible biased translation.

*Identity Preservation.* To preserve the identity/category of the synthesized images, some of the recent works such as [9,46] keep categorical/identity information in discriminator network as an additional task. Some of the others propose to employ a separate classification network which is usually pre-trained [31,57]. In both these cases, the categories/identities are known beforehand and are fixed in number. Thus it is trivial to include such supervision in a GAN framework by training the classifier with real data. However such setup is not feasible in our case as images of new identities to-be-generated are not available to pre-train a classification network (see Section 3.3 for further discussion)

To address the limitation of existing methods of retaining identity/category information of synthesized images, we employ a combination of different set-based supervision approaches for unknown identities to be distinct in the pre-trained embedding space. We keep track of moving averages of same-id features by the momentum-like centroid update rule of center loss [50] and penalize distant same-id samples and close different-id samples by a simplified variant of magnet loss[39] without its sampling process and with only one cluster per identity.

## 3 Adversarial Identity Generation

In this Section, we describe in details the proposed method. Fig. 1 shows the detailed schematic diagram of our method. Specifically, the synthetic image set $x \in \mathcal{S}$ is formed by a graphical engine for the randomly sampled 3DMM, pose and lighting parameters $\alpha$. Then they are translated into more photorealistic domain $G(x)$ through the network $G$ and mapped back to synthetic domain ($G'(G(x))$) through the network $G'$ to retain $x$. Adversarial synthetic and real domain translation of $G$ and $G'$ networks are supervised by the discriminator networks $D_R$ and $D_S$, with an additional adversarial game between $G$ and $G'$ as generator and discriminator respectively. During training, generated identities by 3DMM is preserved with a set-based loss on a pre-trained embedding network $C$. In the following sub-sections, we further describe these components *i.e.* domain adaptation, real-synthetic pair discriminator, and identity preservation.

### 3.1 Unsupervised Domain Adaptation

Given a 3D morphable model (3DMM) [3], we synthesize face images of new identities sampled from its Principal Components Analysis (PCA) coefficients' space with random variation of expression, lighting and pose. Similar to [63], a synthetic input image ($x \in \mathcal{S}$) is mapped to photorealistic domain by a residual network ($G : S \rightarrow \hat{R}$) and mapped back to synthetic domain by a 3DMM fitting network ($G' : \hat{R} \rightarrow \hat{S}$) to complete forward cycle only. To preserve cycle consistency, the resulting image $G'(G(x))$ is encouraged to be the same as input $x$ by a pixel level $L_1$ loss:

$$\mathcal{L}_{cyc} = \mathbb{E}_{x \in \mathcal{S}} \|G'(G(x)) - x\|_1 \tag{1}$$

In order to encourage the resulting images $G(x)$ and $G'(G(x))$ to have similar distribution as real and synthetic domains respectively, those refiner networks are supervised by discriminator networks $D_R$ and $D_S$ with images of the respective domains. The discriminator networks are formed as auto-encoders as in boundary equilibrium GAN (BEGAN) architecture [2] in which the generator and discriminator networks are trained by the following adversarial training formulation:

$$\mathcal{L}_G = \mathbb{E}_{x \in \mathcal{S}} \|G(x) - D_R(G(x))\|_1 \tag{2}$$

$$\mathcal{L}_{G'} = \mathbb{E}_{x \in \mathcal{S}} \|G'(G(x)) - D_S(G'(G(x)))\|_1 \tag{3}$$

$$\mathcal{L}_{D_R} = \mathbb{E}_{x \in \mathcal{S}, y \in \mathcal{R}} \|y - D_R(y)\|_1 - k_t^{D_R} \mathcal{L}_G \tag{4}$$

$$\mathcal{L}_{D_S} = \mathbb{E}_{x \in \mathcal{S}} \|x - D_S(x)\|_1 - k_t^{D_S} \mathcal{L}_{G'} \tag{5}$$

where for each training step $t$ and the network $G$ we update the balancing term with $k_t^{D,G} = k_{t-1}^{D,G} + 0.001(0.5\mathcal{L}_D - \mathcal{L}_G)$. As suggested by [2], this term helps to balance between generator and discriminator and stabilize the training.

### 3.2 Adversarial Pair Matching

Cycle consistency loss ensures bijective transitivity of functions $G$ and $G'$ which means generated image $G(x) \in \hat{R}$ should be transformed back to $x \in \hat{S}$. Convolutional networks are highly under-constrained and they are free to make any unintended changes as long as the cycle consistency is satisfied. Therefore, without additional supervision, it is not guaranteed to achieve the correct mapping that preserves shape, texture, expression, pose and lighting attributes of the face image from domains $S$ to $\hat{R}$ and $\hat{R}$ to $\hat{S}$. This problem is often addressed by introducing pixel-level penalization between input and output of the networks [63,42] which is sub-optimal for domain adaptation as it encourages to stay in the same domain.

To overcome this issue, we propose an additional pair-wise adversarial loss that assign $G'$ network an additional role as a pair-wise discriminator to supervise $G$ network. Given a set of paired synthetic and real images $(\mathcal{P}_\mathcal{S}, \mathcal{P}_\mathcal{R})$, the discriminator loss is computed by BEGAN as follows:

$$\mathcal{L}_{D_P} = \mathbb{E}_{s \in \mathcal{P}_\mathcal{S}, r \in \mathcal{P}_\mathcal{R}} \|s - G'(r)\|_1 - k_t^{D_P} \mathcal{L}_{cyc} \tag{6}$$

While $G'$ network is itself a generator network ($G' : \hat{R} \rightarrow \hat{S}$) with a separate discriminator ($D_S$), we use it as a third pair-matching discriminator to supervise $G$ by means of distribution of paired correspondence of real and synthetic images. Thus while cycle-loss optimizes for biject correspondence, we expect resulting pairs of ($x \in S, G(x) \in \hat{R}$) to have similar correlation distribution as paired training data ($s \in \mathcal{P}_\mathcal{S}, r \in \mathcal{P}_\mathcal{R}$). Fig 2 shows its relation to the previous related arts and comparison to an alternative which is matching aware discriminator with paired inputs for text to image synthesis as suggested by [36]. Please notice that how BEGAN autoencoder architecture is utilized to align the distribution of pair of synthetic and real images with synthetic and generated images.

(a) DC-GAN[35]      (b) BEGAN [2]      (c) Ours      (d) GAN-CLS [36]
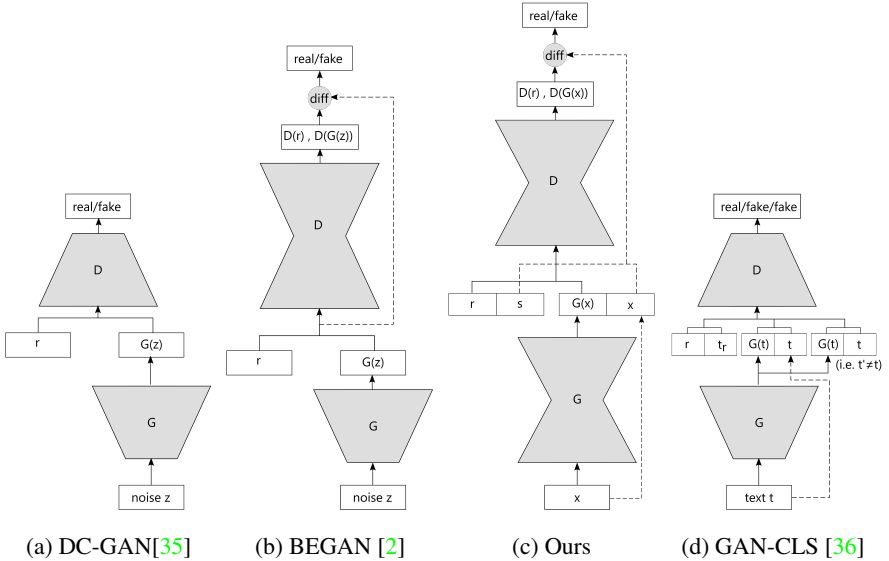
Fig. 2: Comparison of our pair matching method to the related work. (a) In the traditional GAN approach, discriminator module align the distribution of real and synthetic images which is designed as a classification network. (b) BEGAN[2] and many others showed that alignment of error distribution offers more stable training and better results. (c) We propose to utilize this autoencoder approach to align the distribution of pairs to encourage generated image to be a correct transformation to the realistic domain with a game between real and synthetic pairs. (d) An alternative to our method is to introduce wrongly labeled generated images to the discriminator to teach pair-wise matching. [36] used such approach for text to images synthesis.

### 3.3   Identity Preservation

Although identity information is provided by the 3DMM in shape and texture parameters, it may be lost to some extent by virtue of a non-linear transformation. Some studies [57,46] address this issue by employing identity labels of known subjects as additional supervision either with a pre-trained classification network or within the discriminator network. However, we intend to generate images of new identities sampled from 3DMM parameter space and their photorealistic images simply do not exist yet. Furthermore, training a new softmax layer and the rest of the framework simultaneously becomes a chicken-egg problem and results in failed training.

In order to preserve identity on the changing image space, we propose to adapt a set-based approach over a pre-trained face embedding network. We import the idea of pulling same-id samples as well as pushing close samples from different identities in the embedding space such that same-id images are gathered and distinct from other identities regardless of the quality of the images during the training. At the embedding layer of a pre-trained network $C$, generator network ($G$) is supervised by a combination of center [50] and pushing loss [16], which is also a simplified version of Magnet loss [39]

Fig. 3: Quality of 9 images of 3 identities (per row) during the training. Background plot shows the error by the proposed identity preservation layer over the iterations. Notice the changes on the level of fine-details on the faces which is the main motivation of using set-based identity preservation.

formulation which is as following for a given mini-batch (M):

$$\mathcal{L}_C = \mathbb{E}_{x \in \mathcal{S}, i_x \in \mathbb{N}^+} \sum_x^M -log \frac{\exp(\frac{1}{2\sigma^2}\|C(G(x)) - c_{i_x}\|_2^2 - \eta)}{\sum_{j \neq i_x} \exp(\frac{1}{2\sigma^2}\|C(G(x)) - c_j\|_2^2)} \qquad (7)$$

where $i_x$ stands for the identity label of $x$ provided by 3DMM sampling. Margin term $\eta$ is set to 1 and the variance is computed by $\sigma = \frac{\sum_x^M \|C(G(x)) - c_{i_x}\|_2^2}{M-1}$.

While the quality of images is improved during the training, their projection on the embedding space is shifting. In order to adapt to those changes, we update identity centroids ($c_j$) with a momentum of $\beta = 0.95$ when new images of id $j$ is available. Following [50], for a given $x$, moving average of a identity centroid is calculated by $c_j^{t+1} = c_j^t - \beta\delta(i_x = j)(c_j^t - C(G(x)))$ where $\delta(condition) = 1$, if the condition is satisfied and $\delta(condition) = 0$ if not. Centroids ($c_j$) are initialized with zero and after few iterations, they converge to embedding centers and then continue updating to adapt to the changes caused by the simultaneous training of $G$. Fig. 3 shows quality of 9 images of 3 identities over training iterations. Please notice the difference of the images after convergence with the images at the beginning of the training which Softmax layer might converge and fail to supervise for the forthcoming images in later iterations.

**Full Objective**

Overall, the framework is optimized by the following updates simultaneously:

$$\theta_G = \arg\min_{\theta_G} \mathcal{L}_G + \lambda_{cyc}\mathcal{L}_{cyc} + \lambda_C\mathcal{L}_C \qquad (8)$$

$$\theta_{G'} = \arg\min_{\theta_{G'}} \mathcal{L}_{G'} + \lambda_{cyc}\mathcal{L}_{cyc} + \lambda_{D_P}\mathcal{L}_{D_P} \qquad (9)$$

$$\theta_{D_R}, \theta_{D_S} = \arg\min_{\theta_{D_R}, \theta_{D_S}} \mathcal{L}_{D_R} + \mathcal{L}_{D_S} \qquad (10)$$

where $\lambda$ parameters balance the contribution of different modules. The selection of those parameters is discussed in the next section.

## 4    Implementation Details

*Network Architecture:* For the generator networks ($G$ and $G'$), we use a shallow ResNet architecture as in [23] which supplies smooth transition without changing the global structure because of its limited capacity with only 3 residual blocks. In order to benefit from 3DMM images fully, we also add skip connections to the network $G$. Additionally, we add dropout layers after each block in the forward pass with a 0.9 keep rate in order to introduce some noise that could be caused by uncontrolled environmental changes.

We construct the discriminator networks ($D_R$ and $D_S$) as autoencoders trained by boundary equilibrium adversarial learning with Wasserstein distance as proposed by [2]. The classification network $C$, is a shallow FaceNet architecture [41], more specifically we use NN4 network with an input size of $96 \times 96$ where we randomly crop, rotate and flip generated images $G(x)$ which are in size of $108 \times 108$.

*Data:* Our framework needs a large amount of real and synthetic of face images. For real face images, we use CASIA-Web Face Dataset [56] that consists of $\tilde{5}$00K face images of $\tilde{1}$0K individuals.

Please recall that the proposed method trains the $G'$ network as a discriminator ($D_P$) with a small number of paired examples of real and synthetic images. For that, we use a combination of 300W-3D and AFLW2000-3D datasets as our paired training set [64] which consist of 5K real images with their corresponding 3DMM parameter annotations. We render synthetic images by those latent parameters and pair them with matching the real images. This dataset relatively small compared to the ones used by fully supervised transformation GANs (i.e. Amazon Handbag dataset used by [22] contains 137K bag images)

We randomly sample 500K face images of 10K identities as our synthetic data set using Large Scale Face Model (LSFM) [4] and Face Warehouse model for expressions [7]. While shape and texture parameters of new identities are sampled to be under Gaussian distribution of the original model, expression, lighting and pose parameters are sampled with the same Gaussian distribution as synthetic samples of 300W-3D and AFLW2000-3D. For our experiments, we align the faces using MTCNN [60] and centre crop them to the size of $108 \times 108 \times 3$ pixels.

*Training Details:* We train all the components of our framework together from scratch except the classification network $C$ which is pre-trained by using a subset of Oxford VGG Face Dataset [33]. The whole framework takes about 70 hours to converge on a Nvidia GTX 1080TI GPU for 248K iterations with batch size of 16. We start with a learning rate of $8 \times 10^{-5}$ with ADAM solver [25] and halve it at after 128Kth, 192Kth, 224Kth, 240Kth, 244Kth, 246Kth and 247Kth iterations.

As shown in Eqn. 8, 9, $\lambda$ is a balancing factor which controls the contribution of each optimization. We set $\lambda_{cyc} = 0.5$, $\lambda_{D_P} = 0.5$, $\lambda_C = 0.001$ to balance between realism, cycle-consistency, identity preservation and the supervision by the paired data. We also add identity loss ($\mathcal{L}_{id} = \|x - G(x)\|$) as suggested by [63] to regularize the training with a balancing term $\lambda_{id} = 0.1$. During the training, we keep track of moving averages of the network parameters to generate images.

Fig. 4: Random samples from GANFaces dataset. Each row belongs to same identity. Notice the variation in pose, expression and lighting.

As side notes, in our experiments, we observed that it is beneficial to keep non-adversarial signals weak to avoid mode collapse. We also observed that the approach of keeping the history of refined images proposed by [5] breaks adversarial training in our case due to the auto-encoder discriminators.

## 5    Results and Discussions

In this section, we show qualitative and quantitative results of the proposed framework. We also discuss and show the contribution of each module (i.e. $\mathcal{L}_{cyc}$, $D_P$, $C$) with an ablation study in the supplementary materials. For the experiments, we generate 500,000 images of 10,000 different identities with variations on expression, lighting and poses. We name this synthetic dataset **GANFaces**. Please see Fig.4 for random samples from the dataset. The dataset, training code, pre-trained models and face recognition experiments can be viewed at https://github.com/barisgecer/facegan.

### 5.1    Visually Plausible 3DMM Generation

One of the main goals of this work is to generate the face images guided by the attributes of synthetic input images *i.e.* shape, expression, lighting, and poses. We can see from the Fig. 5 that our model is capable of generating photorealistic images preserving the attributes conditioned by the synthetic input images. In the Figure, top row shows the variations of pose and expression on input synthetic faces and the left column shows the input synthetic faces of different identities. And, the rest of the images are the images generated by our model conditioned on the corresponding attributes from top row and left column. We can clearly see that the conditioned attributes are preserved on the
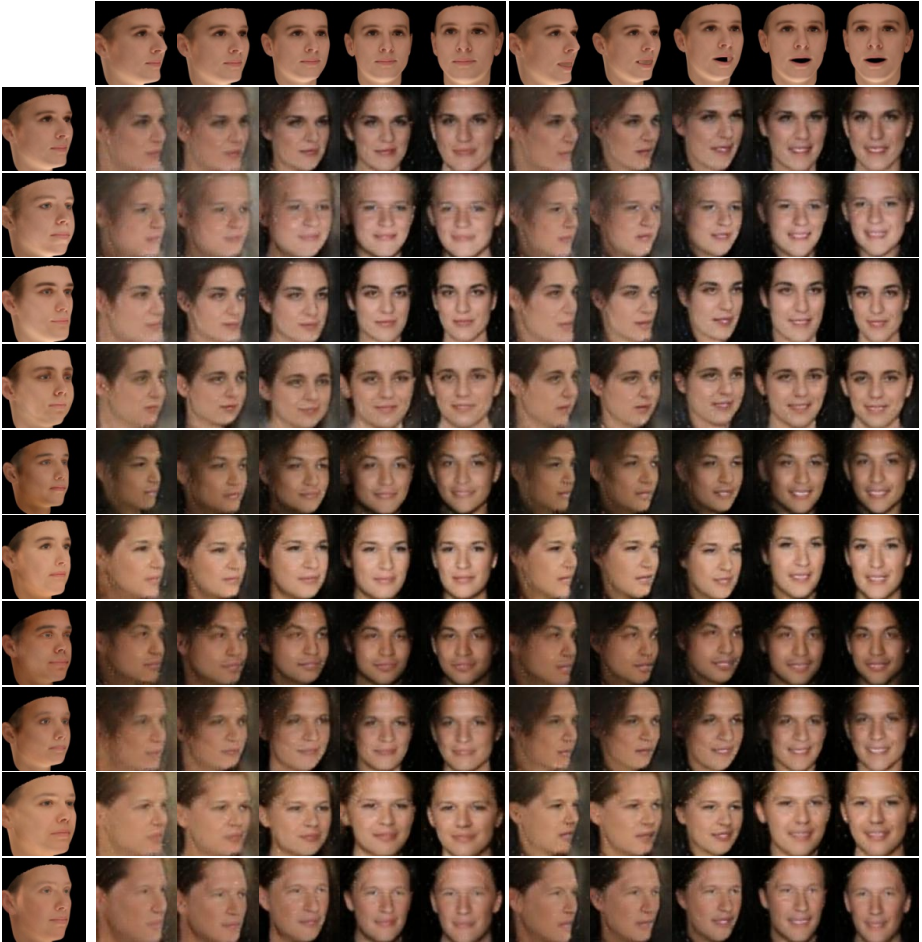
Fig. 5: Images generated by the proposed approach conditioned with identity variation in vertical axis, normalized and mouth open expression in left and right blocks and pose variation in horizontal axis. Images in this figure are not included in the training

images generated by our model. We can also observe that fine-grained attributes such as shapes of chin, nose and eyes are also retained on the images generated by our model. In case of extreme poses, the quality of the image generated by our model becomes less sharp as the CASIA-WebFace dataset, which we used to learn the parameters of discriminator network $D_R$, lacks sufficient number of examples with extreme poses.

## 5.2   The Added Realism and Identity Preservation

In order to show that synthetic images are effectively transformed to the realistic domain with preserving identities, we perform a face verification experiments on GAN-
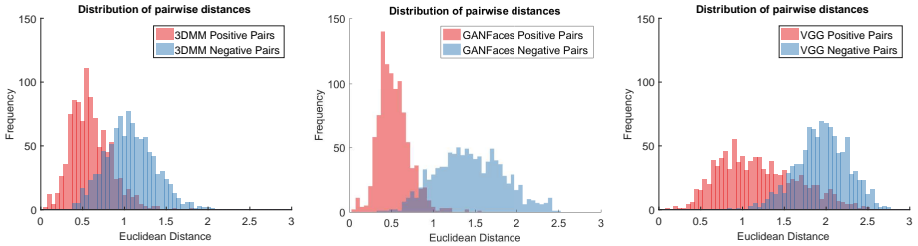
Fig. 6: Distances of 1000 positive and 1000 negative pairs from three different datasets (GANFaces, 3DMM synthetic images, Oxford VGG) embedded on a NN4 network that is trained with CASIA Face dataset

Faces dataset. We took pre-trained face-recognition CNN network, namely FaceNet NN4 architecture [41] trained on CASIA-WebFace [56] to compute the features of the face images. The verification performance of the network on LFW is %95.6 accuracy and %95.5 1-EER which shows that the model is well optimized for in-the-wild face verification. We created 1000 similar (belonging to same identity) and 1000 dissimilar(belonging to different identities) face image pairs from GANFaces. Similarly, we also generated the same number of similar and dis-similar face images pairs from VGG face dataset [33] and the synthetic 3DMM rendered faces dataset. Fig. 6 shows histogram of euclidean distances between similar and dis-similar images measured in the embedding space for the three datasets. The addition of realism and preservation of identities of the GANFaces can be seen from the comparison of its distribution to the 3DMM synthetic dataset's distribution. As the images become more realistic, they become better separable in the pre-trained embedding space. We also observe that the separation of positive and negative pairs of GANFace's faces are better than that of VGG faces pairs. The probable reason of VGG does not having better separation than GANFaces is due to noisy face labels and this is indicated on its original study [33].

### 5.3    Face Recognition with GANFaces dataset

We augmented GANFaces with real face dataset *i.e.* VGG Faces [33] and train VGG19 [43] network and tested performance on two challenging datasets: Labeled Faces in the Wild (LFW) [20] and IJB-A [26]. We restrict ourselves from limited access to full access of real face dataset and train deep network on different combination of real and GAN-Faces. Following [32], we use a pre-trained VGGNet by [43] with 19 layers trained on ImageNet dataset [40] and took these parameters as initial parameters. We train the network with different portion of Oxford VGG Face dataset [33] augmented with the GANFaces dataset. We remove the last layer of deep VGGNet and add two soft-max layers to the previous layer, one for each of the datasets. Learning rate is set to 0.1 for the soft-max layers and 0.01 to the pre-trained layers with ADAM optimizer. Also we halve the gradient coming from GANFaces soft-max. We decrease the learning rate exponentially and train for 80,000 iterations where all of our models are well converged
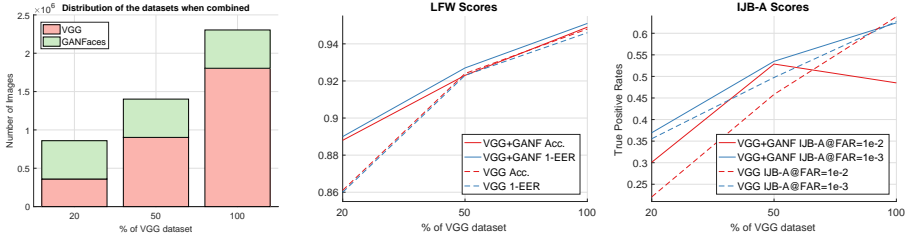
Fig. 7: Face recognition benchmark experiments. (Left) Number of images used from the two datasets in the experiments. Total number of images of VGG Data set is 1.8M since some images were removed from the URL (Middle) Performances on the LFW dataset with and without GANFaces dataset. (Right) True Positive Rates on IJB-A verification task with and without GANFaces dataset.

without overfitting. For a given input size of $108 \times 108$, we randomly crop and flip $96 \times 96$ patches and overall training takes around 9 hours on a NVIDIA 1080TI GPU.

We train 6 models with %20, %50 and %100 of the VGG Face dataset with and without the augmentation of GANFaces. We evaluate the models on LFW and IJB-A datasets and the benchmark scores is improved with the usage of GANFaces dataset even though low resolution images. The contribution of GANFaces increase inversely proportional to the number of images included from VGG dataset which indicates more synthetic images might improve the results even further. Further details can be seen in Fig. 7.

We compare our best model with full VGG dataset and GANFaces to the other state of the art methods. Despite the very low resolution compared to the others, GANFaces was able to improve our baseline to the numbers comparable to the state-of-the-arts. Please note that generative methods such as [32,57], do generation (i.e. pose augmentation and normalization) in the test time where we use only given test images. Together with low resolution, this makes our models more efficient at test time. Given that we only generated 500K images show that the accuracy can be boosted even further by generating more (i.e. 5 times larger from the real set as [32]).

| Method | Real | Synth | Test time Synth | Image size | Acc. (%) | 100% - EER |
|---|---|---|---|---|---|---|
| FaceNet [41] | 200M | - | No | 256×256 | 98.87 | - |
| VGG Face [33] | 2.6M | - | No | 256×256 | 98.95 | 99.13 |
| Masi *et al.* [32] | 495K | 2.4M | Yes | 256×256 | 98.06 | 98.00 |
| Yin *et al.* [55] | 495K | 495K | Yes | 256×256 | 96.42 | - |
| VGG(%100) | 1.8M | - | No | 108×108 | 94.8 | 94.6 |
| VGG(%100) + GANFaces | 1.8M | 500K | No | 108×108 | **94.9** | **95.1** |

Table 1: Comparison with state-of-the-art studies on LFW performances

# 6   Conclusions

In this paper, we propose a novel end-to-end semi-supervised adversarial training framework to generate photorealistic faces of new identities with wide ranges of poses, expressions, and illuminations from 3DMM rendered faces. Our extensive qualitative and quantitative experiments show that the generated images are realistic and identity preserving.

We generated a dataset of 500,000 face images and combined it with a real face image dataset to train a face recognition CNN and improve the performances in recognition and verification tasks. Despite the limited the number of images generated, they were still enough to improve recognition rates. In the future, we plan to generate millions of high resolution images of thousands of new identities to boost the state-of-the-art face recognition.

## Acknowledgments

## References

1. A. Bansal, C. Castillo, R. Ranjan, and R. Chellappa. The do's and don'ts for cnn-based face verification. *ICCVW*, 2017. 1
2. D. Berthelot, T. Schumm, and L. Metz. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017. 2, 6, 7, 9
3. V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194. ACM Press/Addison-Wesley Publishing Co., 1999. 2, 5
4. J. Booth, A. Roussos, S. Zafeiriou, A. Ponniah, and D. Dunaway. A 3d morphable model learnt from 10,000 faces. In *CVPR*, 2016. 9
5. K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. *CVPR*, 2017. 4, 10
6. K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *NIPS*, 2016. 4
7. C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*, 20(3):413–425, 2014. 9
8. Q. Chen and V. Koltun. Photographic image synthesis with cascaded refinement networks. *ICCV*, 2017. 4
9. X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NIPS*, 2016. 5
10. P. Costa, A. Galdran, M. I. Meyer, M. D. Abràmoff, M. Niemeijer, A. M. Mendonça, and A. Campilho. Towards adversarial retinal image synthesis. *arXiv preprint arXiv:1701.08974*, 2017. 2, 4

11. C. Dong, C. C. Loy, K. He, and X. Tang. Image super-resolution using deep convolutional networks. *TPAMI*, 38(2):295–307, 2016. 1

12. A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox. Flownet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 1

13. V. Dumoulin, I. Belghazi, B. Poole, A. Lamb, M. Arjovsky, O. Mastropietro, and A. Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016. 2

14. Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016. 4

15. B. Geçer. *Detection and classification of breast cancer in whole slide histopathology images using deep convolutional networks*. PhD thesis, Bilkent University, 2016. 1

16. B. Gecer, V. Balntas, and T.-K. Kim. Learning deep convolutional embeddings for face representation using joint sample-and set-based supervision. In *ICCVW*, 2017. 3, 7

17. I. Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *NIPS*, 2016. 2, 4

18. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 4

19. K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017. 1

20. G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007. 12

21. X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. *ICCV*, 2017. 4

22. P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 2, 3, 4, 9

23. J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. *ECCV*, 2016. 9

24. T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *ICLR*, 2018. 2

25. D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 9

26. B. F. Klare, B. Klein, E. Taborsky, A. Blanton, J. Cheney, K. Allen, P. Grother, A. Mah, and A. K. Jain. Pushing the frontiers of unconstrained face detection and recognition: Iarpa janus benchmark a. In *CVPR*, 2015. 12

27. C. Lassner, G. Pons-Moll, and P. V. Gehler. A generative model of people in clothing. *ICCV*, 2017. 4

28. J. Li, K. A. Skinner, R. M. Eustice, and M. Johnson-Roberson. Watergan: unsupervised generative network to enable real-time color correction of monocular underwater images. *IEEE Robotics and Automation Letters*, 3(1):387–394, 2018. 4

29. M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NIPS*, 2017. 4

30. W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song. Sphereface: Deep hypersphere embedding for face recognition. In *CVPR*, 2017. 1

31. Y. Lu, Y.-W. Tai, and C.-K. Tang. Conditional cyclegan for attribute guided face image generation. *arXiv preprint arXiv:1705.09966*, 2017. 5

32. I. Masi, A. T. Trn, T. Hassner, J. T. Leksut, and G. Medioni. Do we really need to collect millions of faces for effective face recognition? In *ECCV*, 2016. 4, 12, 13

33. O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *BMVC*, 2015. 1, 9, 12, 13

34. V. M. Patel, R. Gopalan, R. Li, and R. Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE signal processing magazine*, 32(3):53–69, 2015. 2

35. A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 2, 4, 7

36. S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 4, 6, 7

37. S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 1

38. E. Richardson, M. Sela, and R. Kimmel. 3d face reconstruction by learning from synthetic data. In *3D Vision (3DV)*, 2016. 1, 4

39. O. Rippel, M. Paluri, P. Dollar, and L. Bourdev. Metric learning with adaptive density discrimination. *arXiv preprint arXiv:1511.05939*, 2015. 5, 7

40. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 12

41. F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 1, 9, 12, 13

42. A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. *CVPR*, 2017. 2, 4, 6

43. K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 12

44. L. Sixt, B. Wild, and T. Landgraf. Rendergan: Generating realistic labeled data. *arXiv preprint arXiv:1611.01331*, 2016. 2, 4

45. B. Sun and K. Saenko. Subspace distribution alignment for unsupervised domain adaptation. In *BMVC*, 2015. 4

46. L. Tran, X. Yin, and X. Liu. Disentangled representation learning gan for pose-invariant face recognition. In *CVPR*, 2017. 3, 5, 7

47. E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. In *ICCV*, 2015. 4

48. E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. *CVPR*, 2017. 4

49. G. Varol, J. Romero, X. Martin, N. Mahmood, M. Black, I. Laptev, and C. Schmid. Learning from synthetic humans. *CVPR*, 2017. 1, 4

50. Y. Wen, K. Zhang, Z. Li, and Y. Qiao. A discriminative feature learning approach for deep face recognition. In *ECCV*, 2016. 3, 5, 7, 8

51. L. Wolf, Y. Taigman, and A. Polyak. Unsupervised creation of parameterized avatars. *ICCV*, 2017. 4

52. E. Wood, T. Baltrušaitis, L.-P. Morency, P. Robinson, and A. Bulling. Learning an appearance-based gaze estimator from one million synthesised images. In *Proceedings of the Ninth Biennial ACM Symposium on Eye Tracking Research & Applications*, pages 131–138. ACM, 2016. 1, 4

53. C. Xiong, L. Liu, X. Zhao, S. Yan, and T.-K. Kim. Convolutional fusion network for face verification in the wild. *IEEE Transactions on Circuits and Systems for Video Technology*, 26(3):517–528, 2016. 1

54. C. Xiong, X. Zhao, D. Tang, K. Jayashree, S. Yan, and T.-K. Kim. Conditional convolutional neural network for modality-aware face recognition. In *ICCV*, 2015. 1

55. D. Yi, Z. Lei, and S. Z. Li. Towards pose robust face recognition. In *CVPR*, 2013. 13

56. D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 9, 12

57. X. Yin, X. Yu, K. Sohn, X. Liu, and M. Chandraker. Towards large-pose face frontalization in the wild. *ICCV*, 2017. 3, 4, 5, 7, 13

58. S. Yuan, Q. Ye, B. Stenger, S. Jain, and T.-K. Kim. Bighand2. 2m benchmark: Hand pose dataset and state of the art analysis. In *CVPR*, 2017. 1

59. H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017. 4

60. K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503, 2016. 9

61. X. Zhang, Y. Sugano, M. Fritz, and A. Bulling. Appearance-based gaze estimation in the wild. In *CVPR*, 2015. 1, 4

62. Z. Zheng, L. Zheng, and Y. Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. *ICCV*, 2017. 4

63. J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 2, 3, 4, 5, 6, 9

64. X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. In *CVPR*, 2016. 9

# Semi-supervised Adversarial Learning to Generate Photorealistic Face Images of New Identities from 3D Morphable Model: Supplementary Material

Baris Gecer[1], Binod Bhattarai[1], Josef Kittler[2], and Tae-Kyun Kim[1]

[1]Department of Electrical and Electronic Engineering, Imperial College London, UK
[2]Centre for Vision, Speech and Signal Processing, University of Surrey, UK
{b.gecer,b.bhattarai,tk.kim}@imperial.ac.uk,
j.kittler@surrey.ac.uk

## 1 Ablation Study

### 1.1 Quantitative Results

We investigate the contributions of three main components of our framework by an ablation study. Namely, identity preservation module ($\mathcal{L}_C$), adversarial pair matching $\mathcal{L}_{D_P}$ and cycle consistency loss $\mathcal{L}_{cyc}$[1]. We train our framework from scratch in the same way as explained in the paper by removing each of these modules separately (*i.e.* for VGG (%50) version). In table 1, we show the contribution of each module and compare them to the whole framework as a baseline and to the performance of a model trained by only half of the VGG dataset.

| Method | IJB-A Ver. @FAR=0.01 | IJB-A Ver. @FAR=0.001 |
|---|---|---|
| Ours without $\mathcal{L}_C$ | $0.50532 \pm 0.00433$ | $0.17636 \pm 0.00611$ |
| Ours without $\mathcal{L}_{D_P}$ | $0.48034 \pm 0.00402$ | $0.15300 \pm 0.00348$ |
| Ours without $\mathcal{L}_{cyc}$ | $0.49701 \pm 0.00558$ | $0.18341 \pm 0.00670$ |
| VGG (%50) | $0.49751 \pm 0.00484$ | $0.17580 \pm 0.00557$ |
| Ours (VGG(%50)+GANFaces) | $\mathbf{0.53507 \pm 0.00575}$ | $\mathbf{0.18768 \pm 0.00388}$ |

Table 1: Quantitative ablation study. Each of the modules removed from the proposed framework and the performance of generated images are measured on IJB-A verification task

### 1.2 Qualitative Results

Fig. 1 shows visual comparisons between the proposed framework and its versions without each of its components. For the framework and its three variants, we show

---

[1] Here we do not investigate the contribution of the discriminators as their effect is shown by many other studies.

generated images for 12 3DMM input images of 4 different identities with random illumination, pose and expression variations. We evaluate the quality of the images by identity preservation, the visual plausibility and diversity (*i.e.* avoiding mode collapse). Regarding these criteria, our framework clearly generates better images than all of its variants. Without $\mathcal{L}_C$ (Fig.1(c)), namely identity preservation module, the framework *forgets* identity information throughout the network as there is no direct signal to encourage identity preservation. Please notice the identity consistency of our framework (Fig.1(b)) compared to (Fig.1(c)) in the details of faces (*i.e.* shape of nose, eyes, month, eyebrows and their relative distances) or simply by visual gender test. Without adversarial pair matching mechanism $\mathcal{L}_{D_P}$ (Fig.1(d)), we observe local mode collapse across different identities such as shape of nose and eyes in the figure seems to be similar compared to Fig.1(b). This mode collapse is also verified by the quantitative experiments (Table. 1). Cycle consistency loss $\mathcal{L}_{cyc}$ (Fig.1(e)) helps to retain the initial shape given by 3DMM and improves the overall image quality. We also observe noise reduction in the generated images due to the additional supervision provided by all the modules.
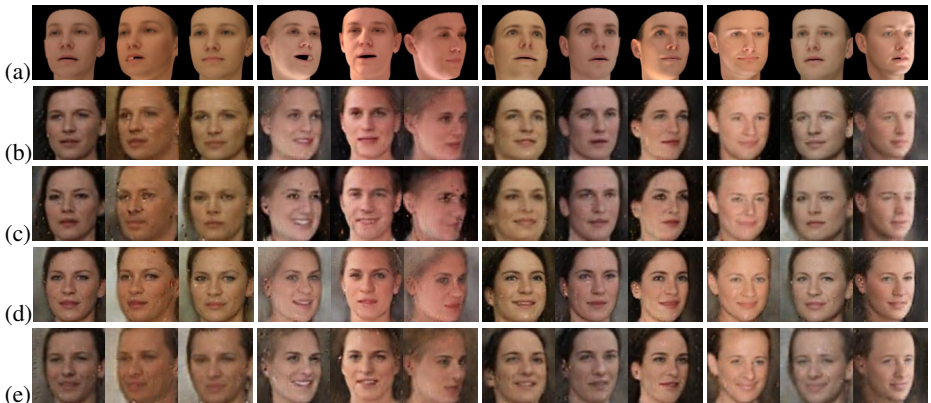


Fig. 1: Columns: divided into blocks of 3 images from the same identity. Rows: (a) 3DMM synthetic images. (b) Generated images by the framework (Ours). (c) Ours without $\mathcal{L}_C$. (d) Ours without $\mathcal{L}_{D_P}$. (e) Ours without $\mathcal{L}_{cyc}$.

## 2  Identity and Illumination Interpolations

In this section, we show the generalization ability of our framework for unseen synthetic identities by interpolating in the identity space of our 3DMM model. Fig . 2 shows how well shapes introduced by the 3DMM is learned so that the transition is smooth and accurate in the photorealistic space. The smooth transition between the two identities with pose variation also shows that the network did not overfit to the given synthetic data and is able to generate more photo-realistic images even without further training.
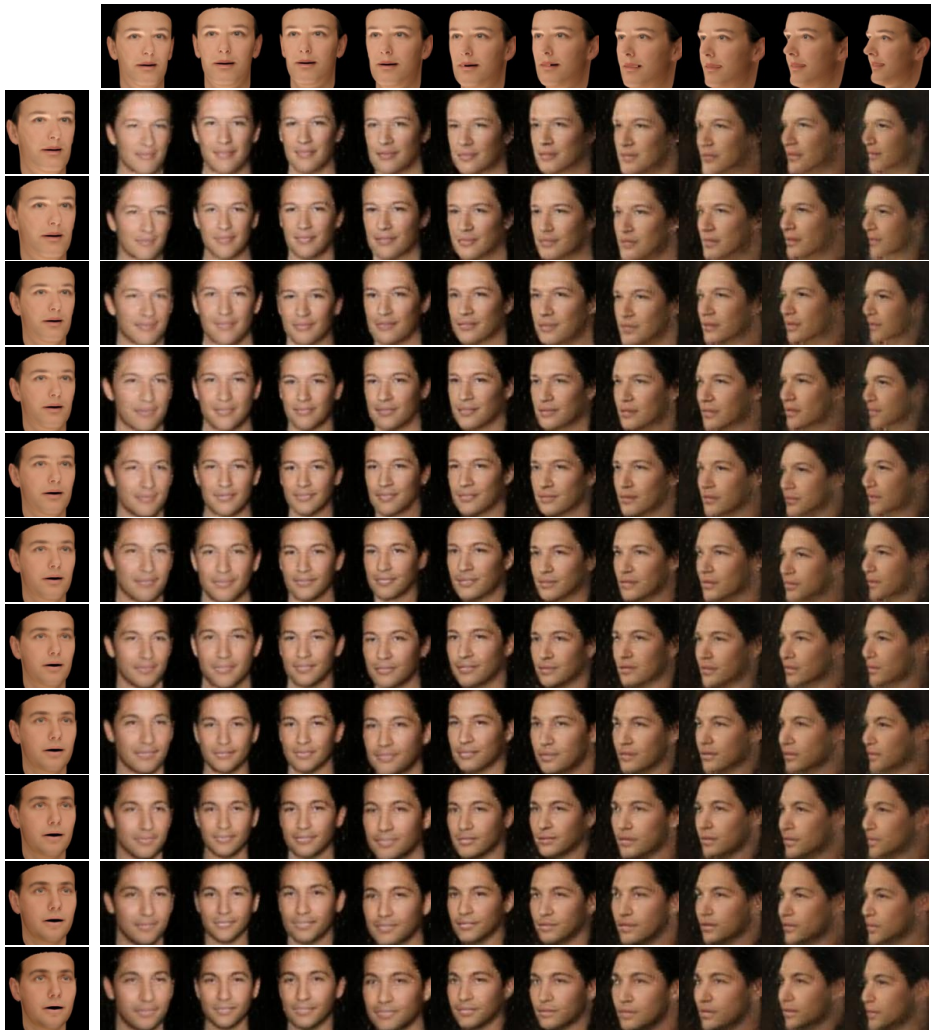
Fig. 2: Identity interpolation between first and second identities of the GANFaces dataset (Fig. 4 first two rows). Interpolation is done in the 3DMM space and projected onto realistic space by our framework. The vertical axis shows the identity interpolation under neutral lighting and expression with pose variation at the horizontal axis. Top-most and left-most 3DMM images indicate the respective identity and pose. Images in this figure are not included in the training.

Figure 3 shows that the framework also learned changes in the illumination strength and able to generate images with a controlled lighting variation.

Fig. 3: Effect of illumination changes to the generated images. Top row contains 3DMM synthetic images and the bottom contains the generated images by the framework given input images as the top row. Extreme lighting conditions result in blurry images as the real training set does not contain images of similar conditions.