

Towards Open-Set Identity Preserving Face Synthesis

Jianmin Bao¹, Dong Chen², Fang Wen², Houqiang Li¹, Gang Hua²

¹University of Science and Technology of China ²Microsoft Research

jmbao@mail.ustc.edu.cn {doch, fangwen, ganghua}@microsoft.com lihq@ustc.edu.cn

Abstract

We propose a framework based on Generative Adversarial Networks to disentangle the identity and attributes of faces, such that we can conveniently recombine different identities and attributes for identity preserving face synthesis in open domains. Previous identity preserving face synthesis is largely confined to synthesizing faces with known identities that are already in the training datasets. To synthesize a face with identity outside the training dataset, our framework requires one input image of that subject to produce an identity vector, and any other input face image to extract an attribute vector capturing, e.g., pose, emotion, illumination, and even the background. We then recombine the identity vector and the attribute vector to synthesize a new face of the subject with the extracted attribute. Our proposed framework does not need to annotate the attributes of the faces at all. It is trained with an asymmetric loss function to better preserve the identity and stabilize the training process. It can also effectively leverage large amount of unlabeled training face images to further improve the fidelity of the synthesized faces for subjects that are not presented in the labeled training face dataset. Our experiments demonstrate the efficacy of the proposed framework. We also demonstrate its usage in a much broader set of applications including face frontalization, face attributes morphing, and face adversarial examples detection.

1. Introduction

Realistic face image synthesis has many real-world applications, such as face super-resolution, frontalization, and morphing, etc.. With the emergence of deep generative models, such as the Generative Adversarial Networks (GAN) [10] and the Variational Auto-encoder [16], we have made tremendous progress in building deep networks for synthesizing realistic faces [34, 19, 33, 19]. However, identity preserving face synthesis remains a challenge, especially when the identity of the face is not presented in the training face images.

Many previous works have attempted to synthesize face

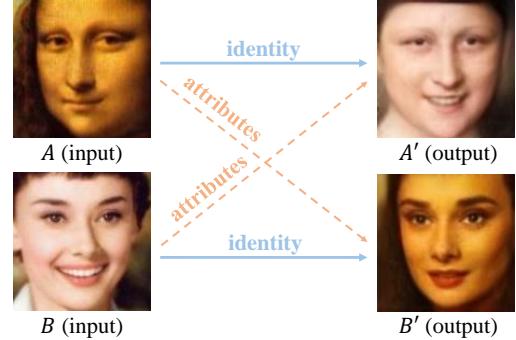


Figure 1. Our method can disentangle identity and attributes from a single face image. With the extracted identity and attributes from two input images A and B , we can generate two new face images A' and B' by recombining the identities and attributes.

images of a specific person. For example, TP-GAN [14] made a stab on synthesizing the frontal view of a face from a single face image. DR-GAN [33] can change the pose of an input face image. However, these methods can only manipulate limited types of attributes, such as poses. Besides, these methods also require full annotation of the attributes for training the models. More recent work, such as CVAE-GAN [4], can produce a variety of attribute changes. Nevertheless, it is not able to synthesize a face with identity outside the training dataset.

In view of the limitation of CVAE-GAN, we make a first stab on the problem of open-set identity preserving face synthesis. Our goal is to be able to synthesize face images of any specific identity, no matter the identity presents in the training dataset or not. To synthesize a face with identity outside the training dataset, we require one input image of that subject to produce an identity vector, and any other input face image to extract an attribute vector capturing, e.g., pose, emotion, illumination, and even the background. We then combine the identity vector and the attribute vector to synthesize a new face of the subject with the extracted attribute.

To this end, we propose a framework based on Generative Adversarial Networks to disentangle identity and attributes given a face image, and recombine different identities and attributes for identity preserving face synthesis.

As shown in Figure 2, our framework has five parts: 1) an identity encoder network I to encode the identities of subjects; 2) an attribute encoder network E to extract attributes of any given face images; 3) a generator network G to synthesize a face image from a combined input of identity and attributes; 4) a classification network C to preserve the identity of the generated face; and 5) a discriminative network D to distinguish real and generated examples. These five parts are trained end-to-end.

In this framework, we propose a new, simple yet elegant way to extract the attributes of an input face image. The proposed framework do not need to annotate the attributes of the faces at all. We use two loss functions: 1) a reconstruction loss of the attribute image, and 2) a KL divergence loss defined on the attribute vector, to enforce the network A to extract the attribute information.

We take full advantage of recent advancement in face recognition, and use the softmax loss on top of the network I to encode the identity into an attribute independent vector representation. Therefore, in order to reconstruct the input, the network A are forced to extract the attribute information. Meanwhile, we add a KL divergence loss to regularize the attribute vector, such that it dose not contain identity information.

Inspired by the CVAE-GAN [4], we adopt a new asymmetric loss function. More specifically, we adopt a cross-entropy loss when training the discriminative network D , and the classification network C , and use a pairwise feature matching loss when update the generative network G . It can better preserve the identity and meanwhile stabilize the training process.

Last but not least, the proposed framework enables us to effectively leverage large amount of unlabeled training face images to further improve the fidelity of the synthesized faces for subjects that are not presented in the labeled training face dataset. These unlabeled data can increase intra-class and inter-class variation of the face distributions, and hence improve the diversity of the synthesized faces. As a result, the generated faces present larger changes in poses and expressions.

In the experiments, we demonstrate the efficacy of the proposed framework. We also show that our model can be applied to other tasks, such as face frontalization, face attributes morphing and face adversarial examples detection.

2. Related Work

We briefly summarize the most related works, ranging from general literature of deep generative models to more specific models on face synthesis.

There have been much recent development on deep generative models, including deterministic generative models [9, 29], Generative Adversarial Network(GAN) [10, 26], Variational Auto-encoder (VAE) [16, 30], and autoregres-

sion networks [18], to list a few. Among them, the most popular one is perhaps the GAN [10, 26].

Many variants of GANs has been proposed to improve the stability of its training process [31, 3, 5]. Meanwhile, there are also many works that have added condition information to the generative network and the discriminative network for conditional image synthesis. The condition information could be a discrete label [24, 7, 4, 19], an reference image [22, 8, 36], and even a text sentence [28, 37, 8].

Due to its abundant useful applications, face is a major focus for image generation. Antipov *et al.* [1] proposed a variant of GANs for face aging. Li *et al.* [21] proposed a method for attribute-driven and Identity-preserving face generation. However, the attribute is only limited to some simple ones. TP-GAN [14] adopted a two-pathway generative network to synthesize frontal faces. Both DR-GAN and TP-GAN obtained impressive results on face frontalization, but they need to explicitly label the frontal faces.

Prior work also explores disentangled representation learning. For example, the DC-IGN [17] used a variational auto-encoder based method to learn the disentangled presentation. However, DC-IGN needs to fix one attribute in one batch training, which also needs explicit annotations of the attributes. Luan *et al.* [33] proposed DR-GAN to learn a generative and discriminative representation, which explicitly disentangle the pose leveraging the pose annotations.

In contrast, this paper proposes an Identity Preserving Generative Adversarial Network framework, which do not require any attribute annotations. This framework disentangles the identity and attributes representations, and then uses different recombination of the representations for identity preserving face synthesis. This disentanglement facilitates us to synthesize faces with identities outside what presented in the training datasets. This addresses a serious limitation of previous deep generative model-based identity preserving face synthesis method [4]. It simply can not generate faces of identities outside the training dataset.

3. Identity Preserving GANs

In this section, we introduce our face synthesis networks. To synthesize a face of any specific identities, our framework requires two input images, *i.e.*, one input image x^s of a certain subject identity, and another input image x^a to extract the attributes, *e.g.*, pose, emotion, illumination, and even background. Our network will synthesize a new face image x' of the subject with the extracted attributes.

As show in Figure 2, our framework is based on the Generative Adversarial Networks. It contains five parts: 1) the identity encoder network I ; 2) the attributes encoder network A ; 3) the generative network G ; 4) the classification network C ; and 5) the discriminative network D . The function of the network I is to extract the identity vector $f_I(x^s)$ from the subject image x^s . The network A is adopted to

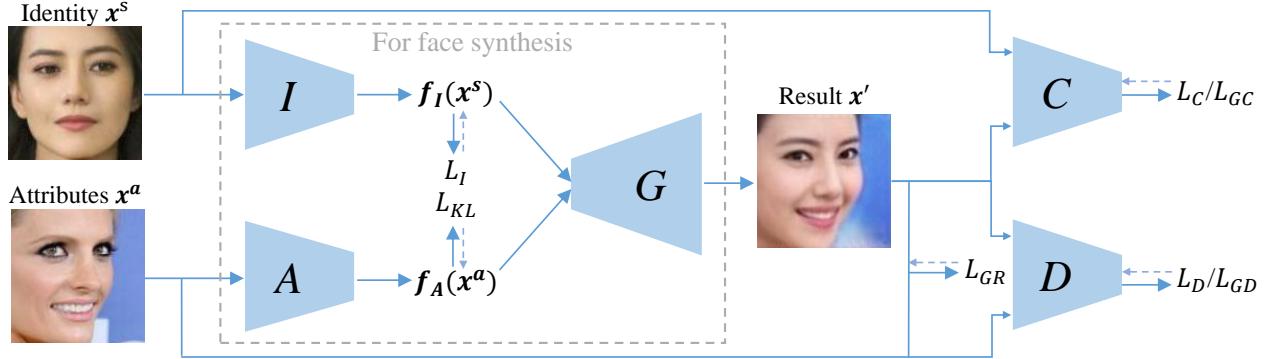


Figure 2. Framework overview: we train a face synthesis network to disentangle identity and attributes from face images and recombine them differently to produce the result x' , which uses the identity of x^s and attributes of x^a . The input/output are drawn with solid lines. And the loss functions are drawn with dashed lines.

extract the attribute vector $f_A(x^a)$ of the attributes image x^a .

The network G generates a new face images x' using the combined identity vector and attribute vector $[f_I(x^s)^T, f_A(x^a)^T]^T$. The network C and the network D are only included in the training phase. The network C is used to preserve the identity by measuring the posterior probability $P(c|x^s)$, where c is the subject identity of x^s . The discriminative network D is to distinguish real and generated examples.

As we do not want to annotate the attributes, how to extract the attribute vector from a face image becomes a challenge. In the following sections, we will first introduce our method on how to disentangle the identity vector and the attribute vector from a face image in Section 3.1. Then, in Section 3.2, we will introduce an asymmetric training method to generate identity-preserving and realistic face images, and to make the training process more stable.

In Section 3.3, in order to further improve the fidelity of the synthesized faces for subjects that are not presented in the labeled training face dataset, we use a unsupervised learning method with large amount of unlabeled training images. Finally, in Section 3.4, we analyze the objective function of the proposed method and provide the implementation details of the training pipeline.

3.1. Disentanglement of Identity and Attributes

In this section, we will introduce the technical details of disentangling the identity vector and the attribute vector using the network I and network A , respectively. In our training data, we only have the annotation of the identity of each face, without any annotation of the attributes information. This is because face images with category annotations are relatively easy to obtain. Many large datasets are publically available, such as the FaceScrub [25], CASIA-WebFace [35] and MS-Celeb-1M [12] datasets. However, the annotation of attributes is often more difficult, and sometimes even impossible, such as the illumination and

the background.

Extracting the identity vector is relatively straightforward. Here, we take full advantage of recent improvement in face recognition. Given a set of face images with the identity annotation $\{x_i^s, c_i\}$, we use the softmax loss for training the network I to perform the face classification task. Therefore, the same individuals have approximately the same feature which can be used as the identity vector. Formally, the loss function of the network I is

$$\mathcal{L}_T = -\mathbb{E}_{x \sim P_r} [\log P(c|x^s)], \quad (1)$$

where $P(c|x^s)$ represents the probability of x^s being with identity c . Then, we use the response of the last pooling layer of I as the identity vector.

In order to train the network A in a fully unsupervised manner to extract the attribute vector. We propose a new, simple yet elegant way to extract the attribute of each face. In the training process, we consider two loss functions: the reconstruction loss and the KL divergence loss.

The reconstruction loss. Here we have two situations: whether the subject image x^s is the same as the attribute image x^a or not. In both cases, we require the result image x' to reconstruct the attribute image x^a , but with different loss weight. Formally, the reconstruction loss is

$$\mathcal{L}_{GR} = \begin{cases} \frac{1}{2} \|x^a - x'\|_2^2 & \text{if } x^s = x^a \\ \frac{\lambda}{2} \|x^a - x'\|_2^2 & \text{otherwise} \end{cases}, \quad (2)$$

where λ is the reconstruction loss weight. Next, we will analyze the reconstruction loss under these two situations.

When the subject image x^s is the same as the attribute image x^a , the output x' is required to be the same as x^s or x^a . Suppose there are various face images of an identity, the identity vector $f_I(x)$ is almost the same for all samples. But the reconstruction using the $f_I(x)$ and $f_A(x)$ of different samples are different. Therefore, the reconstruction loss

will force the attribute encoder network A to learn different attributes representation $\mathbf{f}_A(\mathbf{x})$.

When the subject image \mathbf{x}^s and the attribute image \mathbf{x}^a are different, we do not know exactly what the reconstructed result should look like. But we can expect the reconstruction be approximately similar to the attribute image \mathbf{x}^a , such as the background, the overall illumination, and poses. Therefore, we adopt a raw pixel reconstruction loss with relatively small weight to maintain the attributes. We set $\lambda = 0.1$ in our experiments. As we expected, large λ causes poor result, we further demonstrate this in the supplementary material due to space limit.

The KL divergence loss. To help the attributes encoder network to learn better representations, we also add a KL divergence loss to regularize the attribute vector with an appropriate prior $P(\mathbf{z}) \sim N(0, 1)$. The KL divergence loss will limit the distribution range of the attribute vector, such that it dose not contain much identity information, if at all. For each input face image, the network A outputs the mean μ and covariance of the latent vector. We use the KL divergence loss to reduce the gap between the prior $P(\mathbf{z})$ and the learned distributions, *i.e.*,

$$\mathcal{L}_{KL} = \frac{1}{2} (\mu^T \mu + \sum_{j=1}^J (\exp(\epsilon_j) - \epsilon_j - 1)), \quad (3)$$

where j denotes the j -th element of vector ϵ . Similar to the Variational Auto-encoder [16], we sample the attribute vector using $\mathbf{z} = \mu + \mathbf{r} \odot \exp(\epsilon)$ in the training phase, where $\mathbf{r} \sim N(0, I)$ is a random vector and \odot represents the element-wise multiplication.

3.2. Asymmetric Training for Networks G , C and D

After extracting the identity vector $\mathbf{f}_I(\mathbf{x}^s)$ and attribute vector $\mathbf{f}_A(\mathbf{x}^a)$, we concatenate them in the latent space and feed the combined vector, $\mathbf{z} = [\mathbf{f}_I(\mathbf{x}^s)^T, \mathbf{f}_A(\mathbf{x}^a)^T]^T$ into the network G to synthesize a new face image. In this section, we will introduce our asymmetric training method. This method can generate identity-preserving and realistic face images, which also makes the training process more stable.

Similar to GANs, the generative network G competes in a two-player minimax game with the discriminative network D . The network D tries to distinguish real training data from synthesized data, while the network G tries to fool the network D . Concretely, the network D tries to minimize the loss function

$$\mathcal{L}_D = -\mathbb{E}_{\mathbf{x} \sim P_r}[\log D(\mathbf{x}^a)] - \mathbb{E}_{\mathbf{z} \sim P_z}[\log(1 - D(G(\mathbf{z})))], \quad (4)$$

However, if the network G directly attempts to maximize \mathcal{L}_D as the traditional GAN, the training process of the network G will be unstable. This is because, in practice, the distributions of “real” and “fake” images may not overlap

with each other, especially at the early stage of the training process. Hence, the network D can separate them perfectly. That is, we always have $D(\mathbf{x}^a) \rightarrow 1$ and $D(\mathbf{x}') \rightarrow 0$, as a result, $\mathcal{L}_D \rightarrow 0$. Therefore, when updating the network G , the gradient $\partial \mathcal{L}_D / \partial G \rightarrow 0$. It causes the gradient vanishing problem. Recent works [2, 3] also theoretically analyzed the gradient vanishing problem in training GANs.

To address this problem, inspired by CVAE-GAN [4], we propose to use a pairwise feature matching objective for the generator. To generate realistic face images quality, we match the feature of the network D of real and fake images. Let $\mathbf{f}_D(\mathbf{x})$ denote features on an intermediate layer of the discriminator, then the pairwise feature matching loss is the Euclidean distance between the feature representations, *i.e.*,

$$\mathcal{L}_{GD} = \frac{1}{2} \|\mathbf{f}_D(\mathbf{x}') - \mathbf{f}_D(\mathbf{x}^a)\|_2^2. \quad (5)$$

In our experiment, for simplicity, we choose the input of the last Fully Connected (FC) layer of network D as the feature \mathbf{f}_D .

Meanwhile, the classification network C tries to classify faces of different identities, so the network C tries to minimize the loss function

$$\mathcal{L}_C = -\mathbb{E}_{\mathbf{x} \sim P_r}[\log P(c|\mathbf{x}^s)]. \quad (6)$$

In order to generate identity-preserving face images, we also use pairwise feature matching to encourage \mathbf{x}' and \mathbf{x}^s to have similar feature representations in the network C . Let $\mathbf{f}_C(\mathbf{x})$ denote features produced from an intermediate layer of the classification network C . The feature reconstruction loss is the Euclidean distance between feature representations, *i.e.*,

$$\mathcal{L}_{GC} = \frac{1}{2} \|\mathbf{f}_C(\mathbf{x}') - \mathbf{f}_C(\mathbf{x}^s)\|_2^2. \quad (7)$$

Here, we choose the input of the last FC layer of the network C as the feature for simplicity. We also try to combine features of multiple layers, it only marginally improves the ability to preserve the identity of the network G . The network C and the network I can share the parameters and be initialized by a pretrained face classification network to speed up the convergence.

3.3. Unsupervised Training

Generating face images of identities which do not present in the labeled training dataset remains a challenge. It requires the generative network to cover all intra-person and inter-person variations. Existing publicly available training datasets with labeled identities are often limited by size, and often they do not contain extreme poses or illuminations. In other words, they are not diverse enough.

To solve this problem, we randomly collect about 1 million face images from flicker and Google, and use a face

Net	I	A	G	D	C
Loss	\mathcal{L}_I	$\mathcal{L}_{KL}, \mathcal{L}_{GR}$	$\mathcal{L}_{GR}, \mathcal{L}_{GC}, \mathcal{L}_{GD}$	\mathcal{L}_D	\mathcal{L}_C

Table 1. Networks and their related loss function.

detector to locate the face region. These images have much larger variation, and hence are much more diverse than any existing face recognition datasets. We add these data into the training dataset, and perform an unsupervised training process to help train our generative model to better synthesize face images that do not appear in the training set.

These unlabeled images can be used either as the subject image x^s or the attribute image x^a . When used as the attribute image x^a , the whole training process remains unchanged. When used as the subject image x^s , since it does not have the class label, we ignore the loss function \mathcal{L}_I and \mathcal{L}_C . In other words, the network I and C are fixed, we update the other parts of the end-to-end framework.

These unlabeled data can increase intra-class and inter-class variation of the face distributions, and hence improve the diversity of the synthesized faces. As a result, the generated faces present larger changes in poses and expressions. We will further demonstrate it in the experiments.

3.4. Overall Objective Function

The final synthesis loss function is a sum of all the losses defined above in Equation 1 - 7. Although there are many loss functions, as shown in Table 3.4, each network relates only to one part of the loss function. Therefore, our framework is easy to train and does not need to balance the various loss functions.

In the training phase, we separate each iteration into two steps: one step for the reconstruction process when $x^s = x^a$ and one step for transformation process when $x^s \neq x^a$. The details of the training algorithm is described in Algorithm 1.

4. Experiments

We use a subset of the MS-Celeb-1M [12] dataset, which contains about $5M$ images of $80K$ celebrities for training. For each face image, we first detect the facial region with the JDA face detector [6], and then align and resize it to 128×128 pixels.

For the network I , C , and A , we use the same VGG network [32] structure. Meanwhile, I and C share parameters in the training stage to speed up convergence. For the network G , it is an inverse VGG structure. Its pooling layers are replaced by upsampling layers, and the convolution layers are replaced with deconvolution layers. For the network D , we use the same discriminative network structure as the DCGAN [26]. The batch normalization [15] layer is also applied after each convolution and deconvolution layer. The model is implemented using the deep learning toolbox Torch.

Algorithm 1 Two training process strategy. $\theta_I, \theta_A, \theta_G, \theta_D$, and θ_C are the initial parameters of networks I , A , G , D , and C . $iter \leftarrow 1$.

```

while  $\theta_G$  not converged do
    Sample  $x^s, c$  a batch from the dataset;
    if  $iter \% 2 = 1$  then
        // Training at reconstruction process
         $\lambda \leftarrow 1$ 
         $x^a \leftarrow x^s$ 
    else
        // Training at transformation process
        Sample  $x^a, c$  a batch from the dataset;
         $\lambda \leftarrow 0.1$ 
    end if
     $\mathcal{L}_I \leftarrow -\log(P(c|x^s))$ 
     $\mathcal{L}_C \leftarrow -\log(P(c|x^s))$ 
     $f_I(x^s) \leftarrow I(x^s); f_A(x^a) \leftarrow A(x^a)$ 
     $\mathcal{L}_{KL} \leftarrow KL(f_A(x^a)||P(z))$ 
     $x' \leftarrow G([f_I(x^s)^T, f_A(x^a)^T]^T)$ 
     $\mathcal{L}_D \leftarrow -(log(D(x^a)) + log(1 - D(x')))$ 
     $\mathcal{L}_{GR} \leftarrow \frac{1}{2} \|x^a - x'\|_2^2$ 
     $\mathcal{L}_{GD} \leftarrow \frac{1}{2} \|f_D(x^a) - f_D(x')\|_2^2$ 
     $\mathcal{L}_{GC} \leftarrow \frac{1}{2} \|f_C(x^s) - f_C(x')\|_2^2$ 
     $\theta_I \leftarrow \frac{+}{-} \nabla_{\theta_I}(\mathcal{L}_I)$ 
     $\theta_C \leftarrow \frac{+}{-} \nabla_{\theta_C}(\mathcal{L}_C)$ 
     $\theta_D \leftarrow \frac{+}{-} \nabla_{\theta_D}(\mathcal{L}_D)$ 
     $\theta_G \leftarrow \frac{+}{-} \nabla_{\theta_G}(\lambda \mathcal{L}_{GR} + \mathcal{L}_{GD} + \mathcal{L}_{GC})$ 
     $\theta_E \leftarrow \frac{+}{-} \nabla_{\theta_E}(\lambda \mathcal{L}_{KL} + \lambda \mathcal{L}_G)$ 
     $iter \leftarrow iter + 1$ 
end while

```

4.1. Analysis of the Proposed Framework

In this section, we perform an ablation study of our framework, sweeping loss combinations and different training strategies to understand how each component works in our framework. Both quantitative and qualitative are reported.

We compare five variations of our framework: 1) removing the loss \mathcal{L}_{GD} ; 2) removing the loss \mathcal{L}_{GC} ; 3) training without the transformation training process (denoted as w/o T); 4) training without the unsupervised learning (denoted as w/o U); 5) our best model with all components. The network structure and training strategy remain the same for all settings.

To quantitatively evaluate our framework, we conduct a face identification experiment to compare the performance of each setting. We randomly picked $10K$ people from the MS-Celeb-1M database, each with two photos, one for gallery and one for query. None of these photos are in our training data. For each person, we generate 5 images using the query image and 5 randomly selected attribute images.

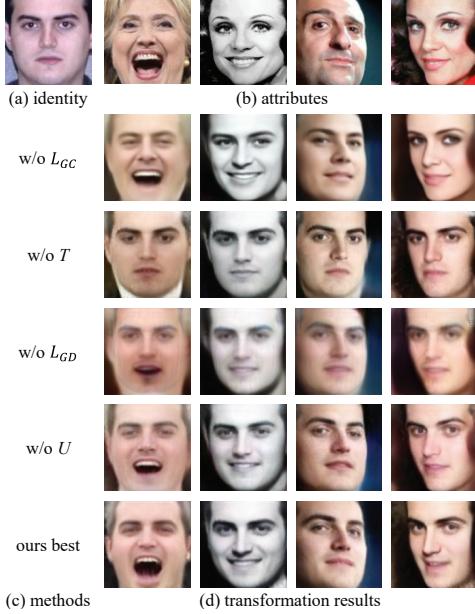


Figure 3. Qualitative comparison between different generative models using different loss combinations and training strategies.

Method	Top-1 acc
original data	87.64
w/o L_{GC}	5.71
w/o T	79.19
w/o L_{GD}	80.52
w/o U	80.24
ours best	81.11

Table 2. Model comparison: Top-1 identification rates (%) using the generated faces.

Then we use the generated image to find the most similar faces in the gallery set, and measure the top-1 accuracy.

In Table 2, we report the face identification top-1 accuracy of each setting. All components can improve the identity preserving capability of the framework. Among them, the L_{GC} contributes the most. Meanwhile, we also measure the top-1 accuracy using the real query image, which is 86.74%. Our generated images achieves comparable results.

Figure 3 illustrates the qualitative results of four variants. We can observe that removing the transformation training process cause the generated results lose attribute details, especially the emotion. Removing the loss L_{GD} will cause the generated image to be blurry. Removing the loss L_{GC} , the generated samples cannot keep the identity information. The results generated by our full model achieves better results. After using unlabeled data for the unsupervised learning, our model can generate images with larger variations. For example, as the bottom left image shown in Figure 3, the mouth can be opened wider.

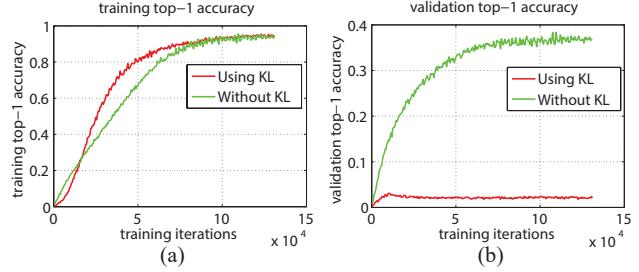


Figure 4. Analysis the KL divergence loss. The attribute vector learned with KL get lower top-1 accuracy.

4.2. Function of KL Divergence Loss

In this section, we will verify whether the KL divergence loss is able to help to remove the identity information in the attribute vector. We first train two models with and without the KL divergence loss, respectively. Then, we conduct the experiment using the FaceScrub [25] dataset. We randomly split these faces into two parts, one for training and the other for validation. For each setting, we use the network A to extract the attribute vector for all the faces in the Facescrub dataset. And then, we use a multilayer perceptron (MLP) to train a classification model to distinguish the face features of different identities in the training set. We also test the top-1 accuracy on the validation set.

The results are presented in Figure 4. We can see that in the validation set, the attribute vector learned using KL loss has lower top-1 accuracy, which means that it contains less identity information. It validates that the KL divergence loss is able to help the network A to remove the identity information in the attribute vector.

4.3. Face Attributes Transformation

This section we show the results of face attributes transformation. The goal of face attributes transformation is to generate an image x' that combines the identity of a input face x^s and the attributes of another input face x^a . To demonstrate that our framework has the ability to generate faces with identities that do not exist in the training dataset. We conduct two experiments: generating face images of identities inside or outside the training dataset, respectively.

Figure 5 presents the face synthesis results of the identities appeared in the training dataset. Our method demonstrated good performance for face synthesis, preserving both on identity and attributes.

Another important feature of our method is that it can synthesize unseen face from the training set. Figure 6 shows the zero-shot identities face synthesis results. Although our model never see these identities, but we can also generate high quality face images which keep the identity and attributes of the given faces.

In Figure 7, we show that our framework can also be used for face frontalization. The results of other methods

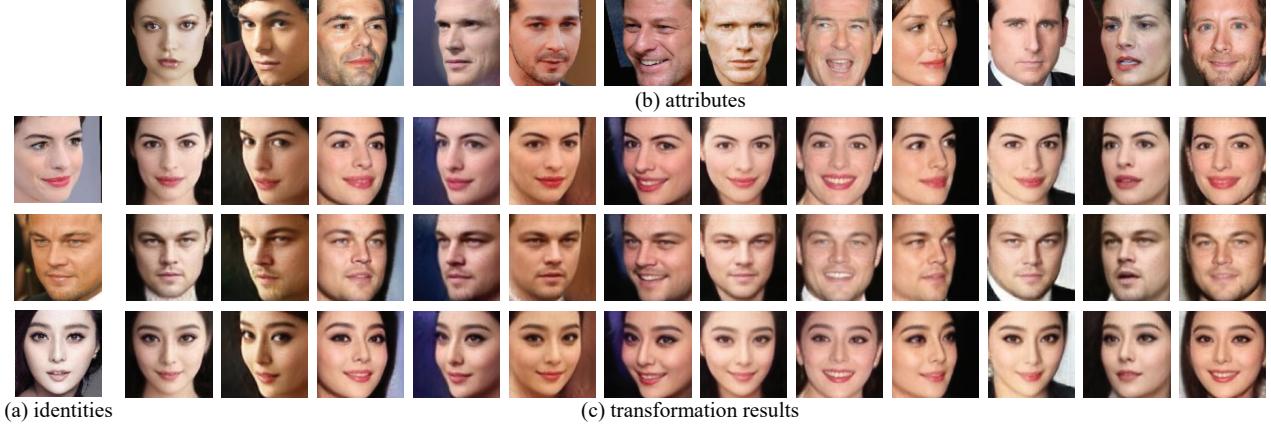


Figure 5. Face synthesis results using the identities appeared in training dataset and random chosen images as attributes.

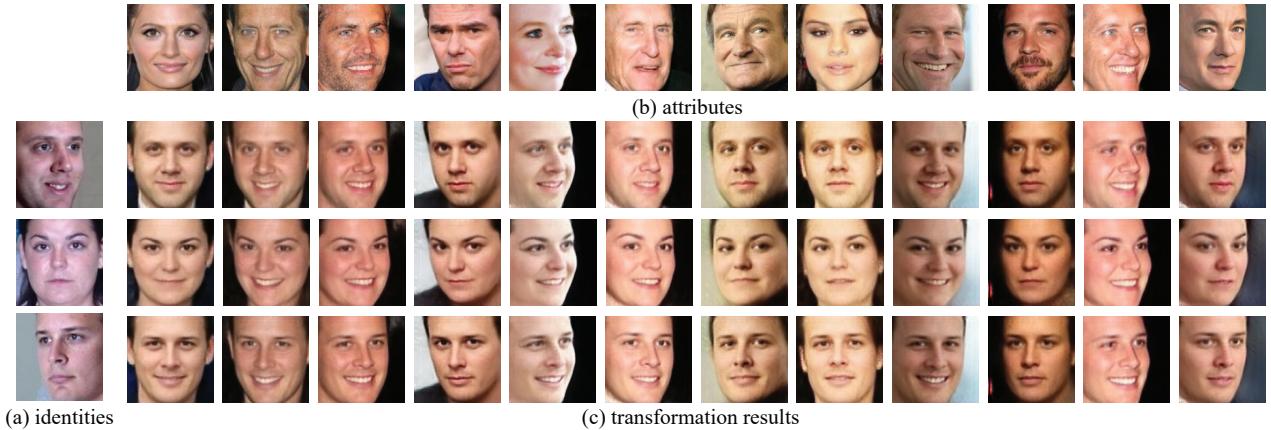


Figure 6. Face synthesis results using the zero-shot identities and random chosen images as attributes.



(a) LFW (b) Ours (c) [14] (d) [38] (e) [13]

Figure 7. Face frontalization results on LFW datasets. Results of other methods are from [14].

are from paper [14]. Although our framework is not trained using pose annotations. With a frontal face as the input image to extract the attributes, our framework can generate identity preserving frontal faces. Comparing with [14], the advantage of our method is that we are able to keep the

lighting and skin color.

4.4. Face Attributes Morphing

In this part, we validate that the attribute in the generated images will continuously change with the latent vector. We call this phenomenon attribute morphing. We test our model on the Multi-PIE [11] dataset. We first select a pair of images x_1 and x_2 , and then extract the attribute vector z_1 and z_2 using the attribute network A . Then, we obtain a series of attribute vectors z by linear interpolation, *i.e.*, $z = \alpha z_1 + (1 - \alpha) z_2, \alpha \in [0, 1]$. Figure 8 presents the results of face attribute morphing. We can gradually change the pose, emotion, or lighting. For more results, please refer to the supplementary material.

4.5. Face Adversarial Examples Detection

Deep network based face verification systems have been widely used in surveillance and access control. However, the existence of adversarial examples make the security or safety of these systems at risk. In this experiment, we will show that our framework can be used for face adversarial examples detection without any modification.

For face verification, given two faces, we first extract the

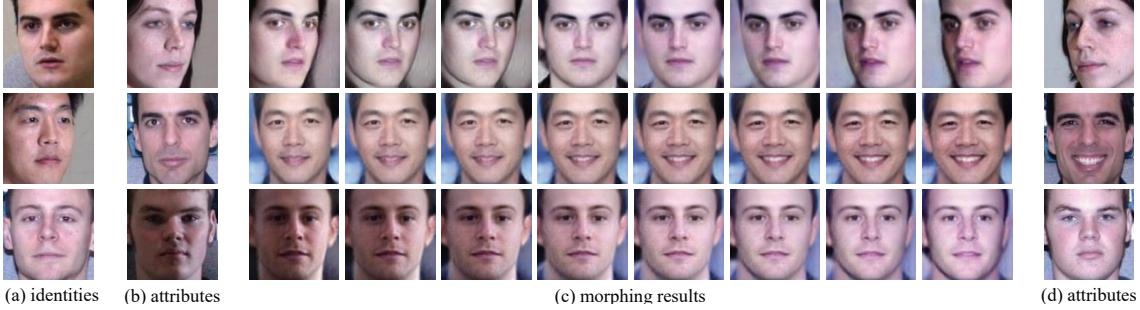


Figure 8. Face morphing results using unseen identities between two attributes. Our framework can gradually change the pose, emotion and lighting.

features for the two faces using a pretrained face classification DNN model. Then, we calculate the distance of the two features, and compare it to a threshold. If the feature distance is smaller than the threshold, they are predicted to have the same identity, and vice versa.

Supposing the two faces x_1 and x_2 have different identities. We can find a imperceptible perturbations r , such that $x_1 + r$ will be regarded as the same person as x_2 using the above face verification system. Here, $x_1 + r$ is the adversarial sample. To find a adversarial sample, we optimize the problem.

$$\begin{aligned} & \min \|r\|_2^2 \\ & \text{s.t. } \|f_C(x_1 + r) - f_C(x_2)\|_2^2 < \theta, \end{aligned} \quad (8)$$

Where f_C is the extracted feature from the pretrained network, and θ is the predefined threshold.

As shown in Figure 9, (a) and (c) are the two inputs x_1 and x_2 , (b) is the adversarial example $x_1 + r$. Since the adversarial examples have similar identity features with others faces, if we reconstruct the image from the feature using the proposed framework, it will generate an image of the other person, as shown in Figure 9. (e). The adversarial example and its reconstruction clearly have different identities. Based on this observation, we can use our generative model to reconstruct the faces, and compare the identity of original faces and reconstruction results to identify the adversarial examples.

We use the LFW [20] to conduct the experiments. For each of 3000 pairs of different identities, we generate two adversarial examples by performing adversarial attack with each other. So, in total we obtain 6000 adversarial examples for each predefined threshold. Here, we choose four different thresholds to conduct the experiments: [0.4, 0.6, 0.8, 1]. At the same time, we have 6000 source images and their reconstruction. Although, we can train another neural network to distinguish the adversarial and source examples [23], the problem is that this neural network can be attacked again. So instead, we choose to use the LBP [27] feature. We extract LBP features from the input and its reconstruction image and then concatenate them. Finally, we train a linear SVM to conduct the binary classification. The

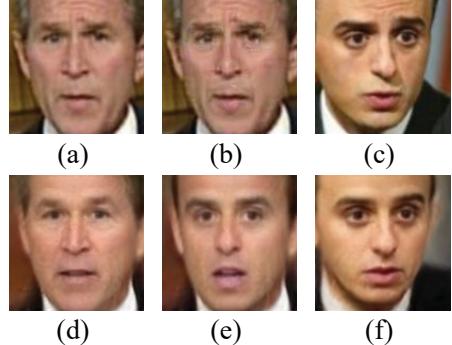


Figure 9. Adversarial examples detection in face verification system. (a) is the source image, (b) is the adversarial example which aims to attack face image (c). (d), (e) and (f) are the reconstruction results of them from our framework. We can clearly observe that although the adversarial example have similar appearance with the source image, but their reconstruction results have different appearance.

threshold	1.0	0.8	0.6	0.4
acc	76.73%	82.58%	87.18%	92.41%

Table 3. Results of adversarial examples detection at different thresholds.

results are shown in Table 3, we can achieve 92.41% accuracy if we require the feature distance to be less than 0.4.

5. Conclusion

In this paper, we propose an Open-Set Identity Preserving Generative Adversarial Network framework for disentangling the identity and attributes of faces, and synthesizing faces from the recombined identity and attributes. The framework shows superior performance in generating realistic and identity preserving face images, even for identities outside the training dataset. Our experiments demonstrate that our framework can also be applied in other tasks: such as face image frontalization, face attributes morphing, and adversarial examples detection in face verification system. In the future work we hope to explore whether our framework can be applied to other datasets, like birds, flowers and rendered chairs.

References

- [1] G. Antipov, M. Baccouche, and J.-L. Dugelay. Face aging with conditional generative adversarial networks. *arXiv preprint arXiv:1702.01983*, 2017. 2
- [2] M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. In *NIPS 2016 Workshop on Adversarial Training. In review for ICLR*, volume 2016, 2017. 4
- [3] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017. 2, 4
- [4] J. Bao, D. Chen, F. Wen, H. Li, and G. Hua. Cvae-gan: Fine-grained image generation through asymmetric training. *arXiv preprint arXiv:1703.10155*, 2017. 1, 2, 4
- [5] D. Berthelot, T. Schumm, and L. Metz. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017. 2
- [6] D. Chen, S. Ren, Y. Wei, X. Cao, and J. Sun. Joint cascade face detection and alignment. In *European Conference on Computer Vision*, pages 109–122. Springer, 2014. 5
- [7] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2172–2180, 2016. 2
- [8] H. Dong, S. Yu, C. Wu, and Y. Guo. Semantic image synthesis via adversarial learning. *arXiv preprint arXiv:1707.06873*, 2017. 2
- [9] A. Dosovitskiy, J. Springenberg, M. Tatarchenko, and T. Brox. Learning to generate chairs, tables and cars with convolutional networks. 2016. 2
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014. 1, 2
- [11] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker. Multi-pie. *Image and Vision Computing*, 28(5):807–813, 2010. 7
- [12] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *European Conference on Computer Vision*, pages 87–102. Springer, 2016. 3, 5
- [13] T. Hassner, S. Harel, E. Paz, and R. Enbar. Effective face frontalization in unconstrained images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4295–4304, 2015. 7
- [14] R. Huang, S. Zhang, T. Li, and R. He. Beyond face rotation: Global and local perception gan for photorealistic and identity preserving frontal view synthesis. *arXiv preprint arXiv:1704.04086*, 2017. 1, 2, 7
- [15] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 5
- [16] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 1, 2, 4
- [17] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum. Deep convolutional inverse graphics network. In *Advances in Neural Information Processing Systems*, pages 2539–2547, 2015. 2
- [18] H. Larochelle and I. Murray. The neural autoregressive distribution estimator. In *AISTATS*, volume 1, page 2, 2011. 2
- [19] A. B. L. Larsen, S. K. Sønderby, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015. 1, 2
- [20] E. Learned-Miller, G. B. Huang, A. RoyChowdhury, H. Li, and G. Hua. Labeled faces in the wild: A survey. In *Advances in Face Detection and Facial Image Analysis*, pages 189–248. Springer, 2016. 8
- [21] M. Li, W. Zuo, and D. Zhang. Convolutional network for attribute-driven and identity-preserving human face generation. *arXiv preprint arXiv:1608.06434*, 2016. 2
- [22] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. *arXiv preprint arXiv:1703.00848*, 2017. 2
- [23] J. H. Metzen, T. Genewein, V. Fischer, and B. Bischoff. On detecting adversarial perturbations. *arXiv preprint arXiv:1702.04267*, 2017. 8
- [24] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2
- [25] H.-W. Ng and S. Winkler. A data-driven approach to cleaning large face datasets. In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 343–347. IEEE, 2014. 3, 6
- [26] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 2, 5
- [27] M. A. Rahim, M. S. Azam, N. Hossain, and M. R. Islam. Face recognition using local binary patterns (lbp). *Global Journal of Computer Science and Technology*, 2013. 8
- [28] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396*, 2016. 2
- [29] S. E. Reed, Y. Zhang, Y. Zhang, and H. Lee. Deep visual analogy-making. In *Advances in neural information processing systems*, pages 1252–1260, 2015. 2
- [30] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014. 2
- [31] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. *arXiv preprint arXiv:1606.03498*, 2016. 2
- [32] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [33] L. Tran, X. Yin, and X. Liu. Disentangled representation learning gan for pose-invariant face recognition. In *CVPR*, volume 4, page 7, 2017. 1, 2
- [34] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. *arXiv preprint arXiv:1512.00570*, 2015. 1
- [35] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 3

- [36] D. Yoo, N. Kim, S. Park, A. S. Paek, and I. S. Kweon. Pixel-level domain transfer. In *European Conference on Computer Vision*, pages 517–532. Springer, 2016. [2](#)
- [37] H. Zhang, T. Xu, H. Li, S. Zhang, X. Huang, X. Wang, and D. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *arXiv preprint arXiv:1612.03242*, 2016. [2](#)
- [38] X. Zhu, Z. Lei, J. Yan, D. Yi, and S. Z. Li. High-fidelity pose and expression normalization for face recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 787–796, 2015. [7](#)