

Geometry Guided Adversarial Facial Expression Synthesis

Lingxiao Song^{1,2} Zhihe Lu^{1,3} Ran He^{1,2,3} Zhenan Sun^{1,2} Tieniu Tan^{1,2,3}

¹National Laboratory of Pattern Recognition, CASIA

²Center for Research on Intelligent Perception and Computing, CASIA

³Center for Excellence in Brain Science and Intelligence Technology, CAS

Abstract

Facial expression synthesis has drawn much attention in the field of computer graphics and pattern recognition. It has been widely used in face animation and recognition. However, it is still challenging due to the high-level semantic presence of large and non-linear face geometry variations. This paper proposes a Geometry-Guided Generative Adversarial Network (G2-GAN) for photo-realistic and identity-preserving facial expression synthesis. We employ facial geometry (fiducial points) as a controllable condition to guide facial texture synthesis with specific expression. A pair of generative adversarial subnetworks are jointly trained towards opposite tasks: expression removal and expression synthesis. The paired networks form a mapping cycle between neutral expression and arbitrary expressions, which also facilitate other applications such as face transfer and expression invariant face recognition. Experimental results show that our method can generate compelling perceptual results on various facial expression synthesis databases. An expression invariant face recognition experiment is also performed to further show the advantages of our proposed method.

1. Introduction

Facial expression synthesis is a classical graphics problem where the goal is to generate face images with specific expression for specified human subject. It has drawn much attention in the field of computer graphics, computer vision and pattern recognition. Synthesizing photo-realistic facial expression images has been of great value for both academic and industrial communities, and has been widely applied in facial animations, face editing, face data augmentation and face recognition. During the last two decades, many facial expression synthesis methods have been proposed, which can be roughly divided into two categories. The first category mainly resorts to computer graphics technique to directly warp input faces to target expressions [39, 35, 37] or re-use sample patches of existing im-

ages [23], while the other aims to build generative models to synthesize images with predefined attributes [30, 6].

For the first category, a lot of research efforts have been devoted to finding correspondence between existing facial textures and target images. Earlier approaches usually generate new expressions by creating fully textured 3D facial models [26, 1], warping face images via feature correspondence [31] and optical flow [35, 36], or compositing face patches from an existing expression dataset [23, 14]. Particularly, Yeh et al. [37] propose to learn the optical flow with a variational autoencoder. Although this kind of methods can usually produce realistic images with high resolution, their elaborated yet complex processes often result in expensive computation.

The representative methods in the second category are deep generative models that have recently obtained impressive results for image synthesis applications [40, 13, 12]. However, images generated by such methods sometimes lack fine details and tend to be blurry or of low-resolution. Targeted expressional attributes are usually encoded in a latent feature space, where certain directions are aligned with semantic properties. Therefore, these methods can provide better flexibility in semantical-level image generation, but it is hard to take fine-grain control of the synthesized images, e.g., widen the smile or narrow the eyes.

In this paper, a deep architecture (G2-GAN) is proposed to synthesize photo-realistic and identity-preserving facial images while keep operation-friendly. A human face is often assumed to contain geometry and texture information [21] in computer vision, and both geometry and texture attributes can be used to facilitate face recognition and expression classification [19]. Inspired by the face geometry information in active appearance models (AAM), we employ face geometry to control the expression synthesis process. Face geometry is defined via a set of feature points, and is transformed to an image (heat map) and fed to G2-GAN as a control condition. Fig. 1 is the pipeline of our approach. We generate facial expression images conditioned on both the input face images and geometry attributes. Particularly, expression generating and removal are simultane-

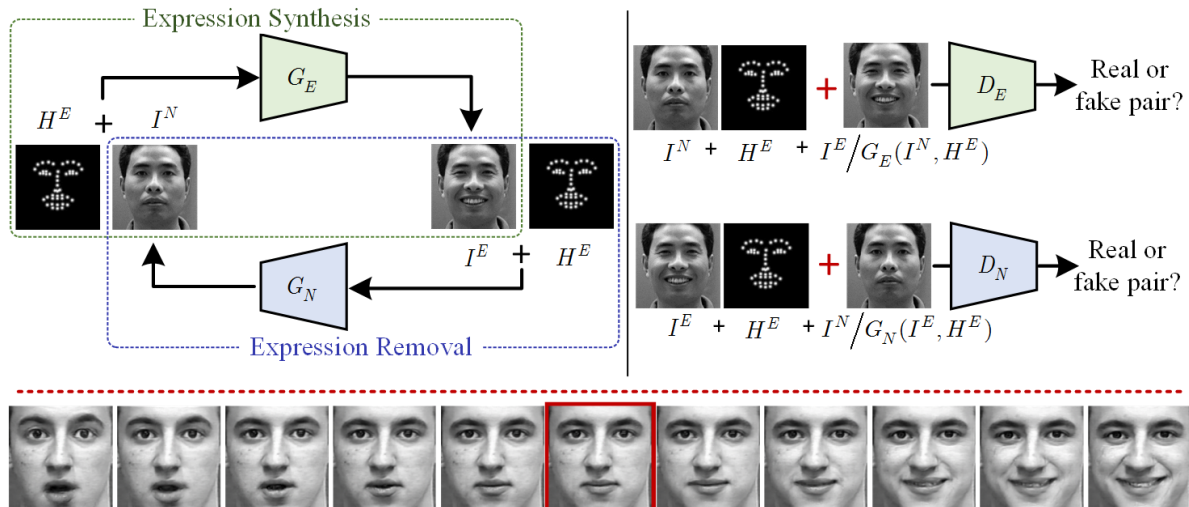


Figure 1. The proposed geometry-guided facial expression generation framework. Face geometry is fed into generators as the condition to guide the processes of expression synthesis and expression removal. In the bottom, we show some examples generated from the same real face image (the center one marked with red box).

ously considered in our method, constructing a full mapping cycle of expression editing. Then, expression transfer can be performed between arbitrary expressions and subjects. Extensive experiments on two facial expression databases demonstrate the superiority of the proposed facial expression synthesis framework.

The main contributions are summarized as follows,

- We propose a novel geometry-guided GAN architecture for facial expression synthesis. It can generate photo-realistic images in different expressions from a single image, where target expressions can be easily controlled by various facial geometry inputs.
- We employ a pair of GANs to simultaneously perform two opposite tasks: removing expression and synthesizing expression. By combining these two models, our method can be used in many applications such as facial expression transfer and cross-expression recognition.
- We utilize an individual-specific shape model to operate facial geometry, which gives consideration to individual differences when perform expression synthesis. Based on this model, facial expression transfer and interpolation can be easily conducted.
- Extensive experiments on two facial expression databases demonstrate that the proposed method can synthesize photo-realistic and identity-preserving expression images.

2. Related Works

Facial expression synthesis (or editing) is an important task in face editing. In this section, we briefly review some recent advances in facial expression synthesis and its related generative adversarial networks (GAN).

2.1. Expression synthesis

As mentioned above, existing expression synthesis methods can be categorized to two classes according to the way of manipulating pixels.

Methods in the first category address this problem either with 2D/3D image warping [1, 8], flow mapping [36, 37] or image reordering [17, 35], most of which are morph-based or example-based. For instance, [1] estimates 3D shape from a neutral face, and synthesizes facial expression by 3D rendering. Bolkart et al. [2] propose a groupwise multilinear correspondence optimization to iteratively refine the correspondence between different 3D faces. In [8], an image-based warping strategy is introduced to perform automatic face reenactment, with the facial identity preserving being considered. Thies et al. [32] track expressions based on a statistical facial prior, and then achieve real-time facial reenactment by using deformation transfer in a low-dimensional expression space. Particularly, Olszewski et al. [24] employ a generative adversarial framework to refine 3D texture correspondences and infer details such as wrinkles and inner mouth region. Many works attempt to utilize the optical flow map to perform image warping. In [36], 3D faces of different expressions are constructed, and expression flow is computed by projecting the difference between 3D shapes back to 2D. Recently, neural networks based

methods [7, 37] have been presented to manipulate expression flow maps. It is difficult for those warping-based methods to recover unseen facial components, e.g., skin wrinkles and inner mouth area, or synthesize realistic images for new faces.

Example-based methods edit faces by re-using image patches or reordering image samples of a training set, which can synthesize expected expressions as well as generate unseen faces. [23] composites facial patches from a large dataset to synthesize face images with desired expressions. In [14], expression is mapped to a new face by matching images with similar pose and expression from a database of the target person. Li et al. [17] hallucinate face videos by retrieving frames via a carefully-designed expression similarity metric from an existing expression database. Yang et al. [35] reorder input face frames using Dynamic Time Warping, and then apply an additional expression warping to get more realistic results.

The other kinds of methods use generative models to deal with facial expression synthesis. In [30], a deep belief net is used to convert high-level descriptions of facial attributes into realistic face images. Reed et al. [28] propose a higher-order Boltzmann machine to model interaction among multiple groups of hidden units, and each unit group encodes distinct variation factors such as pose, morphology and expression in face images. In [5], a regularization term is embedded in an autoencoder to disentangle the variations between identity and expression in an unsupervised manner. Li et al. [18] build a convolutional neural network to generate facial images with the given source input images and reference attribute labels. Shu et al. [29] learn a face-specific disentangled representation of intrinsic face properties via GAN, and generate new faces by changing the latent representations. Recently, Ding et al. [6] propose ExprGAN which can synthesize facial expression with controllable intensity, and an expression controller network is proposed to learn expression code. ExprGAN is the most similar work to ours as far as we know. However, ExprGAN generates images conditioned on expression labels and intensity values, while we employ the face geometry as control condition which is not limited to certain expression styles.

2.2. Generative adversarial networks

Our work is also related to the generative adversarial networks (GAN) [9], which provides a simple yet efficient way to train powerful model via the min-max two player game between generator and discriminator. Many modified architectures of GAN have been proposed to deal with different tasks. For example, CGAN[22] introduces a conditional version of GAN to guide image synthesis process via adding supervised information to both generator and discriminator. CycleGAN [40], DualGAN [38] and DiscoGAN [15] share the same idea of employing a cycle structure to handle the

unpaired image-to-image translation problem. GAN and its variants have achieved great success in numerous image-generating-related tasks such as image synthesis [27], image super-resolution [16], image style transfer [40, 13] and face synthesis [12]. Motivated by this, we develop our facial expression synthesis framework based on GAN, aiming at generating photo-realistic images with high-quality local details.

3. Methods

In this section, we present a novel framework for the facial expression synthesis problem based on generative adversarial networks. We first describe the geometry guided facial expression synthesis in detail, and then propose geometry manipulation methods for face transfer and expression interpolation.

3.1. Geometry Guided Facial Expression Synthesis

The outstanding performance of GAN in fitting data distribution has significantly promoted many computer vision applications such as image style transfer [40, 13]. Motivated by its remarkable success, we employ GAN to perform the facial expression synthesis.

Only limited expression styles are supported by existing deep learning-based facial expression synthesis methods, which are usually semantic properties such as smile and angry. Many works can transform a neutral face to a smile face, but can hardly control how strong the smile is. Even though one can construct an intensity-sensitive model by using training data with emotion intensity annotations, many expressions are still difficult to encode with the limited semantic properties. For example, it is hard to describe “a lopsided grin with one eye open” using normal semantic properties. To address this problem, we employ the face geometry to guide the generation.

As in AAM, face geometry is defined via a set of fiducial points [21]. Heatmap is used to encode the locations of these facial fiducial points, which has been widely used in human pose estimation [33] and face alignment [11]. The heatmap provides a per-pixel likelihood for fiducial point locations. Given the heatmaps of target facial expressions and frontal-looking faces without expression (in the following we term it as expressionless faces), new face images (expressed faces) are synthesized accordingly.

As illustrated in Fig. 1, a pair of generators $G_E : (I^N, H^E) \rightarrow I^E$ and $G_N : (I^E, H^E) \rightarrow I^N$ are introduced, in which I^N is an expressionless face, I^E is an expressed face and H^E is the heatmap corresponding to I^E . Associated with these two generators, two discriminators D_E and D_N are involved, aiming to distinguish between real triplets (I, H, I') and generated triplets $(I, H, G(I))$ correspondingly. I and I' are images of expressionless and expressed faces, or vice versa.

It is worth noting that H^E plays different roles in these two face editing models, i.e., control measure in expression synthesis and auxiliary annotation in expression removal. In the expression synthesis process, H^E is used to specify the target expression so that G^E can transform neutral expression I^N into desired expression. As for the expression removal process, H^E is in charge of indicating the state of I^E so as to facilitate the recovering of I^N .

Adversarial Loss. Generators and discriminators are trained alternatively towards adversarial goals, following the pioneering work of [9]. Since the proposed face editing models generate results conditioned on the input face images and heatmaps, we apply GAN in conditional setting as [22, 13]. The adversarial losses for generator and discriminator are shown in Eq. 1 and Eq. 2 respectively.

$$L_{G-adv} = -\mathbb{E}_{I,H \sim P(I,H)} \log D(I, H, G(I, H)) \quad (1)$$

$$L_{D-adv} = \mathbb{E}_{I,H,I' \sim P(I,H,I')} \log(1 - D(I, H, I')) + \mathbb{E}_{I,H \sim P(I,H)} \log D(I, H, G(I, H)) \quad (2)$$

Pixel Loss. The generator is tasked to not only fool the discriminator, but also synthesize images similar to the target ground-truths as far as possible. The pixel-wise loss L_{pixel} enforces the transformed face image to have a small distance with the ground-truth in the raw-pixel space. L_{pixel} takes the form:

$$L_{pixel} = \mathbb{E}_{I,H,I' \sim P(I,H,I')} \|I' - G(I, H)\|_1, \quad (3)$$

where we use L1 distance to encourage less blurring output. (I, H, I') is one of the combination of (I^N, H^E, I^E) and (I^E, H^E, I^N) depending on the generators.

Cycle-Consistency Loss. The generators G_E and G_N construct a full mapping cycle between neutral expression faces and expressed faces. If we transform a face image from neutral expression to angry and then transform it back to neutral expression, the same face image should be obtained in the ideal situation. Therefore, we introduce an extra cycle consistency loss L_{cyc} to guarantee the consistency between source images and the reconstructed images, e.g., I^N vs. $G_N(G_E(I^N, H^E), H^E)$ and I^E vs. $G_E(G_N(I^E, H^E), H^E)$. L_{cyc} is calculated as

$$L_{cyc} = \mathbb{E}_{I,H \sim P(I,H)} \|I - G'(G(I, H))\|_1, \quad (4)$$

where G' is the opposite generator to G . In our case, if G is used to transform neutral expression into expression specified by the face geometry heatmap H , then G' is used to recover the neutral expression with the assistance of H .

Identity Preserving Loss. A fundamental principle of facial expression editing is that face identity should be preserved after expression synthesis as well as removal. Thus,

an identity-preserving term is adopted in our framework to enforce identity consistency:

$$L_{identity} = \mathbb{E}_{I,H \sim P(I,H)} \|F(I) - F(G(I, H))\|_1, \quad (5)$$

where F is a feature extractor for face recognition. We employ the model-B of the Light CNN [34] as our feature extraction network, which includes 9 convolution layers, 4 max-pooling layers and one fully-connected layer. The Light CNN is pre-trained as a classifier to distinguish between tens of thousands of identities, so it has ability to capture the most prominent feature for face identity discrimination. Therefore, we can leverage this loss to enforce preserving face identity through the face editing processes.

To sum up, the final full objective for generators G_N, G_E is a weighted sum of all the losses defined above: L_{G-adv} to remove the modality gap between real and generated samples, L_{pixel} to force pixel-wise correctness, L_{cyc} to guarantee cycle consistency of the reconstructed image and source image, and $L_{identity}$ to preserve identity characteristic through mapping process.

$$L_G = L_{G-adv} + \alpha_1 L_{pixel} + \alpha_2 L_{cyc} + \alpha_3 L_{identity} \quad (6)$$

where α_1, α_2 and α_3 are loss weight coefficients.

3.2. Facial Geometry Manipulation

As mentioned above, geometric positions of a set of fiducial points are employed to guide facial expression editing in our framework. Face geometry is largely affected by facial expression, and is a useful cue for expression recognition [19]. Its usage provides a more intuitive yet efficient way for specifying target facial expression. This is because face geometry can not only visually represent the locations and shapes of facial organs, but also be adjusted continuously to obtain expressions with different intensities.

Human faces have unique physiological structure characteristics, resulting in strong correlation between the locations of fiducial points. Hence, the variance of facial geometry should be constrained to avoid unreasonable settings, e.g., eyebrows under the eyes, square-shapes eyes or nose. Taking the prior knowledge of faces' distribution into account, a parametric shape model is built to serve as a geometry generator.

We adopt a method similar to [21] to learn a basic shape model from labelled training images. Firstly, faces are normalized to the same scale and rotated to horizontal according to the locations of two eyes. Then, Principal Component Analysis (PCA) is applied to get a basic shape model of the locations for K fiducial points

$$s(p) = s_0 + Sp \quad (7)$$

where $s, s_0 \in R^{2K \times 1}$, $S \in R^{2K \times N}$, $p \in R^{N \times 1}$. The base shape s_0 is the mean shape of all the training images and

columns of S are the N eigenvectors corresponding to the N largest eigenvalues. Different facial geometries can be obtained by changing the value of shape parameters p .

However, facial geometry is not only correlated with facial expression, but also related to face identity to a great extent. The facial geometry varies with different individuals even under the same expression. For example, the distance between eyes and the length of nose depend largely on face identity rather than expression. Considering these individual differences, we propose an individual-specific shape model based on Eq. 7, which can be derived by replacing the mean shape s_0 with the neutral shape s_0^I of different individuals. The individual-specific shape model is given by

$$s^I(p) = s_0^I + Sp \quad (8)$$

where s_0^I accounts for variation relate to identity, while p accounts for changes caused by facial expression.

Facial Expression Transfer. The proposed framework can be easily applied in facial expression transfer. Given two expressed faces I^A and I^B with detected facial landmarks s^A, s^B . The expression removal model is firstly employed to recover expressionless faces as

$$I_0^A = G_N(I^A, s^A), I_0^B = G_N(I^B, s^B) \quad (9)$$

where I_0^A, I_0^B denote the neutral expression faces of I^A, I^B respectively. Therefore the neutral shapes s_0^A, s_0^B can be acquired via facial landmark detection.

Then, the shape parameters are derived by solving the following least squares regression problem.

$$\begin{aligned} p^A &= \arg \min_p \|s^A - s_0^A - Sp\|^2 \\ p^B &= \arg \min_p \|s^B - s_0^B - Sp\|^2 \end{aligned} \quad (10)$$

We change shape parameters so as to get transferred locations of fiducial points.

$$\begin{aligned} s^{AB} &= s_0^A + Sp^B \\ s^{BA} &= s_0^B + Sp^A \end{aligned} \quad (11)$$

Heatmaps are transformed according to these transferred shapes, and concatenated with corresponding expressionless faces as inputs for expression synthesis. Finally, results of facial expression transfer can be obtained by using our expression synthesis model as Eq. 12.

$$I^{AB} = G_E(I_0^A, s^{AB}), I^{BA} = G_E(I_0^B, s^{BA}) \quad (12)$$

Facial Expression Synthesis and Interpolation. As mentioned above, our method is able to synthesize different expressions from a single image. The simple requirement is to prepare a neutral expression face image and shape parameters for target expression. Benefitting from the proposed

expression removal model, neutral expression face is not hard to access. The shape parameters for specific expression can be learnt via the basic shape model (see in Eq. 7) from annotated training dataset. Once the values of shape parameters are associated with certain semantic properties, such as fear and surprise, we can use them to synthesize unseen facial expressions with desired semantic types. Besides, facial expression interpolation can be conducted by linearly adjusting the value of shape parameters.

4. Experiments

In this section, we evaluate the proposed approach on two commonly used facial expression databases. The databases and testing protocols are introduced firstly. Then, the implementation details are presented. Finally, we provide experiments with qualitative and quantitative results for single-image editing, face transfer, expression interpolation and expression-invariant face recognition.

4.1. Datasets and Protocols

The CK+ database [20]. CK+ database includes 593 sequences from 123 subjects, in which seven kinds of emotions are labeled. The first frame is always neutral while the last frame has the peak expression. In each expression video sequence, the first frame is selected as the neutral expression, while last half frames are used as target expression. Training and testing subsets are divided based on identity, with 100 for training and 23 for testing. Locations of 68 fiducial points of each frame are provided, and we use them to create heatmaps for experiments. Because almost all of the videos in CK+ database are grayscale, grayscale images are used in our experiment.

The Oulu-CASIA NIR-VIS facial expression database [4]. Videos of 80 subjects with six typical expressions and three different illumination conditions are captured in both NIR and VIS imaging systems in this database. Only images captured by a VIS camera within strong illumination condition are used in our facial expression editing experiments. Similar to the CK+ database, we take the first frame and images belong the last half of each sequence to make training pairs. The Oulu-CASIA database includes two parts captured among different ethnic groups at different time, where P001 to P050 are Finnish people and the rest P051 to P080 are Chinese people. We find that these two parts differ a lot in illuminations and face structures. Hence, we select training data over these two parts. Finally we get a training subset of 60 subjects that consists of 37 Finns and 13 Chinese, and a testing subset with 13 Finns and 7 Chinese accordingly. We use the 68 fiducial points detected by [3] to create heatmaps.

4.2. Implementation Details

Image pre-processing. All the face images are normalized by the similarity transformation using the locations of two eyes, and then cropped to 144×144 size, of which 128×128 sized sub images are selected by random cropping in training and center cropping in testing. In training stage, we also perform random flipping of the input images to encourage generalization performance. The heatmap is a multi-channel image with the same size as input face image, where value of each pixel is the likelihood for fiducial point location. 2D Gaussian convolution is applied on each channel to smooth the heatmap. All the pixel values are normalized into range of $[0,1]$, including face images and heatmaps.

Network architecture. We adapt our architecture from [13]. The generators take the architecture of U-Net, which is an encoder-decoder with skip connections between mirrored layers in the encoder and decoder stacks. For discriminator networks, the frequently-used PatchGAN model is employed.

We train different models for each dataset with a batch size of 5 and an onset learning rate of 10^{-4} . In all our experiments, hyper-parameters are set empirically to balance the importance of different losses. The trade-off parameter α_1 for pixel loss is set to 10, α_2 for cycle consistency loss is set to 5. α_3 for identity-preserving loss is set to 0.1 in the beginning, and is gradually increased to 0.5 along with the training process.

4.3. Experimental Results

4.3.1 Facial Expression Editing

For this experiment, given testing image triplets (I^N, I^E, H^E) , we conduct expression synthesis on (I^N, H^E) and expression removal on (I^E, H^E) simultaneously. Some visual examples are shown in Fig. 2 and Fig. 3. The first two rows display original expressionless faces and original expressed faces, and the next two rows are results of expression removal and expression synthesis respectively. We can see that the proposed G2-GAN is capable of generating compelling identity-preserving faces for desired expression in both testing datasets. Since the images in the CK+ database have higher resolution than those in the Oulu-CASIA database, results for the CK+ database contain better low-level image quality such as skin wrinkles. Noting that we can synthesize satisfactory mouth region with even teeth textures, without needing to involve extra manipulations such as recovering mouth area by retrieving similar frames from a pre-trained database.

In order to measure the correctness of transformed images, we adopt PSNR (peak signal to noise ratio, dB) and SSIM (structural similarity index) for quantitative metric, where PSNR is calculated on the luminance channel and



Figure 2. Results of CK+ database for facial expression synthesis and removal. From top to bottom, input expressionless images (true I^N), input expressed images (true I^E), expression removal results (fake I^N) and expression synthesis results (fake I^E).

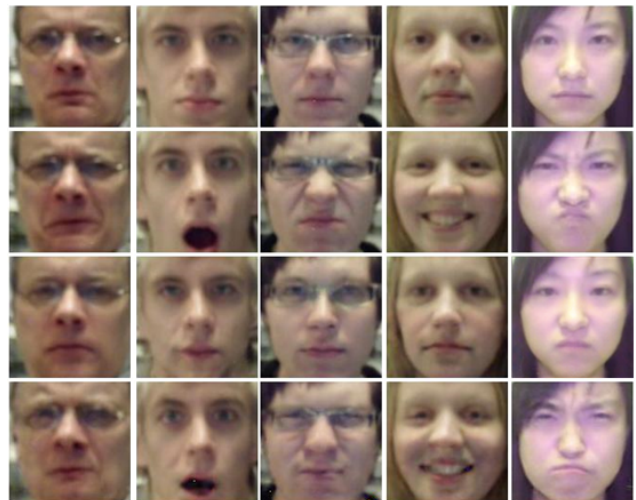


Figure 3. Results of Oulu-CASIA database for facial expression synthesis and removal. Images are arranged by the same order as Fig. 2.

SSIM is calculated on three channels of RGB respectively. Tab. 1 reports quantitative results of the proposed approach under different settings. Both the cycle consistency loss and the identity preserving loss contribute to improve performances, and the best result is acquired by combining them together.

Table 1. Quantitative results for expression synthesis and expression removal on CK+ and Oulu-CASIA databases.

Dataset	Configuration	Expression Removal		Expression Synthesis	
		SSIM	PSNR	SSIM	PSNR
CK+	w/o $L_{cyc}, L_{identity}$	0.726	22.655	0.756	23.903
	w/o L_{cyc}	0.728	22.828	0.754	24.061
	w/o $L_{identity}$	0.724	22.516	0.765	24.335
	G2-GAN	0.728	22.968	0.767	24.420
Oulu-CASIA	w/o $L_{cyc}, L_{identity}$	0.902	25.202	0.908	26.206
	w/o L_{cyc}	0.903	25.270	0.914	26.337
	w/o $L_{identity}$	0.904	25.519	0.916	26.677
	G2-GAN	0.910	25.810	0.914	26.588



Figure 4. Results of CK+ database for facial expression transfer. There are three images for each subject in each example. From the left to right, the input images, results of expression removal, results of facial expression transfer.

4.3.2 Facial Expression Transfer

In this part, we demonstrate our model’s ability to transfer the expression of different faces. The procedures for facial expression transfer are introduced in Sec. 3.2.

Fig. 4 and Fig. 5 show some example results. The facial expressions are transferred between two subjects in an identity consistent way. Besides, identity-irrelevant face attributes, e.g., eyeglasses and hairs, are perfectly preserved. Individual differences are considered in facial expression transfer, resulting in various local deformations for different subjects. For example, when different people keep the same expression of smile, more obvious changes can be discovered for people with larger mouths.

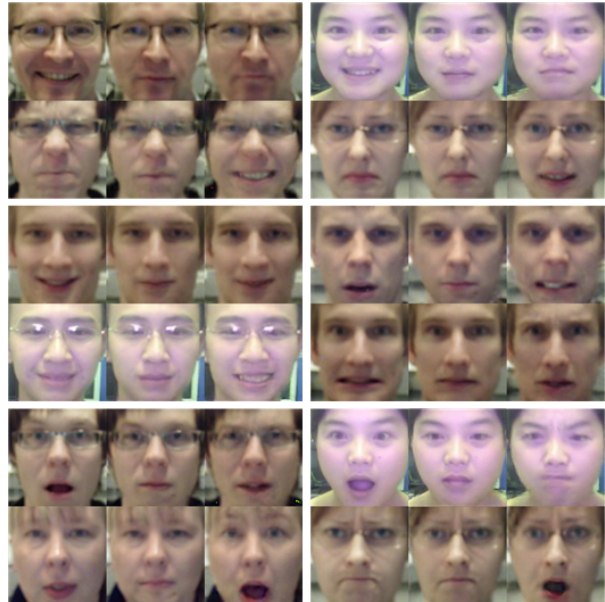


Figure 5. Results of Oulu-CASIA database for facial expression transfer. Images are arranged by the same order as Fig. 4.

4.3.3 Facial Expression Interpolation

Interpolation for unseen expression is conducted in this experiment to demonstrate our model’s capability to synthesize expressions with different intensities. It is worth noting that there is no ground-truth in this experiment, and the locations of the fiducial points are obtained from a pre-trained shape dictionary as described in Sec. 3.2.

The generated images are shown in Fig. 6 and Fig. 7, in which each row contains a new type of expressions with different intensities. G2-GAN successfully transforms the input faces to new unseen expressions with fine details. Especially for the results on the CK+ database, the changes of facial textures caused by expression change are well captured such as glabellar wrinkles under expressions of anger and disgust, chin wrinkles when mouth shut and brows lifting when scared. This validates that the proposed G2-GAN’s adjustability in generating multiple face expressions, not

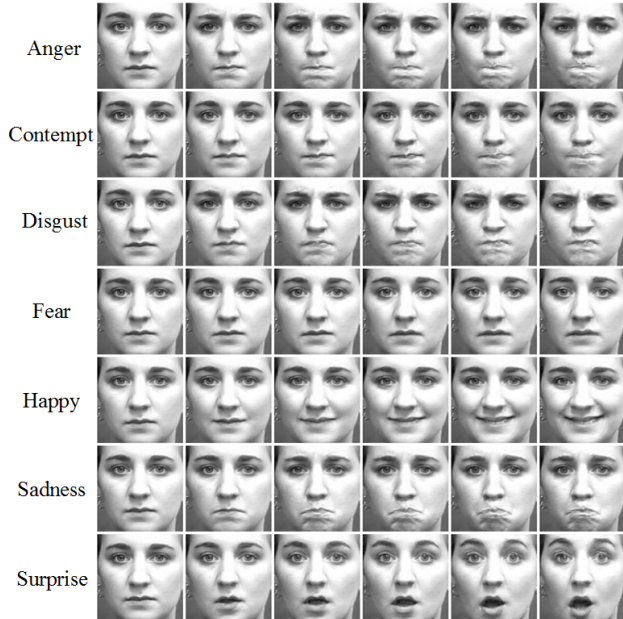


Figure 6. Results of CK+ database for facial expression interpolation. Images in the left-most column are the source images, and the remainder are synthesized results. Each row shows a different expression with ascending intensity from left to right. Seven expression styles are shown corresponding to the annotated expression classes in CK+ database.

limited in pre-determined categories. Besides, these results also demonstrate the operation-friendliness of our method, as we can easily synthesize expressions of desired intensities. An interesting phenomenon is that our model can distinguish the deformations of the mouth caused by happiness and surprise, and the teeth are only generated when synthesizing a smile expression.

4.3.4 Expression-Invariant Face Recognition

In this subsection, we apply G2-GAN in expression-invariant face recognition. The expression removal model is employed as a normalization module in face recognition, which transforms faces into neutral expression. Face verification is taken in both the CK+ dataset and the Oulu-CASIA dataset. The gallery set is selected from the first frame of each video sequences, with only one image for each subject. The probe set is made up of all the rest images in testing set. Two released face recognition models are tested, including the VGG-FACE [25] and the Light CNN [34]. The Rank-1 identification rate, true accept rates at 1% and 0.1% (TAR@FAR=1%, TAR@FAR=0.1%) are taken as evaluation metrics. In order to validate the effectiveness of L_{cyc} and $L_{identity}$, we report the results of removal each one of them respectively.

Results for the expression-invariant face recognition ex-

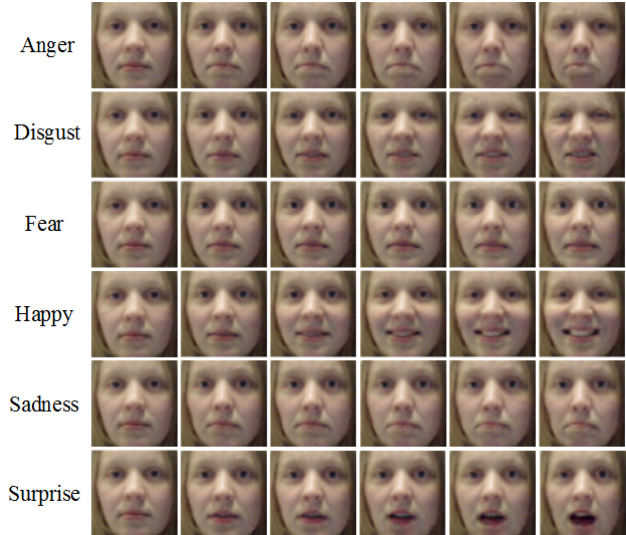


Figure 7. Results of Oulu-CASIA database for facial expression interpolation. Images are arranged by the same order as Fig. 6. Six expression styles are corresponding to the annotated expression classes in Oulu-CASIA database.

periment are presented in Tab. 2. Benefiting from the powerful representation ability of deep learning methods, VGG-FACE and Light CNN obtain high performances on the original images. However, results can be further improved by introducing our expression removal module, especially for a lower FAR. Both the cycle consistency loss and the identity preserving loss facilitate to improve the recognition performance according to results of w/o L_{cyc} , w/o $L_{identity}$, and the basic setting w/o $L_{cyc}, L_{identity}$. Besides, slight drops occur when we do not use $L_{identity}$ comparing with the results of original images, suggesting the necessity of $L_{identity}$ in face editing when the face identity is expected to be preserved.

5. Conclusions

This paper has developed a geometry-guided adversarial framework for facial expression synthesis. Facial geometry has been employed to guide photo-realistic face synthesis as well as to provide an operation friendly solution for specifying target expression. Besides, a pair of facial editing subnetworks are trained together towards two opposite tasks: expression removal and expression synthesis, forming a mapping cycle between expressionless and expressed faces. By combining these two subnetworks, our method can be used in many face related applications including facial expression transfer and expression-invariant face recognition. Moreover, we have proposed an individual-specific shape model for operating the facial geometry, in which individual differences are considered. Extensive experimental results demonstrate the effectiveness of the proposed

Table 2. Results for expression-invariant face recognition on CK+ and Oulu-CASIA databases. Images in the probe set are processed by our expression removal model firstly, and then fed to face recognition models. We conduct face verification on the transformed probe set and the original gallery set. Results of the ‘original’ configuration are obtained by directly testing on the non-transformed gallery set as well as probe set.

Dataset	Configuration	VGG-Face			Light CNN		
		Rank-1	FAR=1%	FAR=0.1%	Rank-1	FAR=1%	FAR=0.1%
CK+	original	96.41	92.13	88.11	100.00	97.01	93.33
	w/o $L_{cyc}, L_{identity}$	96.15	93.33	84.94	98.63	96.83	87.77
	w/o L_{cyc}	96.15	94.27	87.25	100.00	97.60	94.87
	w/o $L_{identity}$	96.41	92.90	84.86	99.23	97.43	89.65
	G2-GAN	97.26	96.15	92.22	100.00	97.69	94.95
Oulu-CASIA	original	97.68	94.63	90.91	99.92	95.35	89.02
	w/o $L_{cyc}, L_{identity}$	96.95	94.99	90.46	99.52	95.95	90.10
	w/o L_{cyc}	97.56	95.80	92.90	99.92	97.60	91.67
	w/o $L_{identity}$	96.59	95.35	90.02	99.84	97.04	89.78
	G2-GAN	97.84	96.19	93.19	99.88	97.80	93.31

method for facial expression synthesis.

References

- [1] V. Blanz, C. Basso, T. Poggio, and T. Vetter. Reanimating faces in images and video. In *Computer Graphics Forum*, pages 641–650, 2003. **1, 2**
- [2] T. Bolkart and S. Wuhler. A groupwise multilinear correspondence optimization for 3d faces. In *ICCV*, 2015. **2**
- [3] A. Bulat and G. Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *ICCV*, 2017. **5**
- [4] J. Chen, D. Yi, J. Yang, G. Zhao, S. Z. Li, and M. Pietikainen. Learning mappings for face synthesis from near infrared to visual light images. In *CVPR*, 2009. **5, 10**
- [5] B. Cheung, J. A. Livezey, A. K. Bansal, and B. A. Olshausen. Discovering hidden factors of variation in deep networks. In *ICLRW*, 2015. **3**
- [6] H. Ding, K. Sricharan, and R. Chellappa. Exprgan: Facial expression editing with controllable expression intensity. In *AAAI*, 2018. **1, 3, 10, 11**
- [7] Y. Ganin, D. Kononenko, D. Sungatullina, and V. Lempitsky. Deepwarp: Photorealistic image resynthesis for gaze manipulation. In *ECCV*, 2016. **2**
- [8] P. Garrido, L. Valgaerts, O. Rehmsen, T. Thormahlen, P. Perez, and C. Theobalt. Automatic face reenactment. In *CVPR*, 2014. **2**
- [9] 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. **3, 4**
- [10] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker. Multi-pie. In *FG*, pages 1–8, 2008. **11**
- [11] L. Huang, Y. Yang, Y. Deng, and Y. Yu. Densebox: Unifying landmark localization with end to end object detection. *arXiv preprint arXiv:1509.04874*, 2015. **3**
- [12] 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. In *ICCV*, 2017. **1, 3**
- [13] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. **1, 3, 4, 6**
- [14] I. Kemelmacher-Shlizerman, E. Shechtman, E. Shechtman, and S. M. Seitz. Being john malkovich. In *ECCV*, 2010. **1, 3**
- [15] T. Kim, M. Cha, H. Kim, J. Lee, and J. Kim. Learning to discover cross-domain relations with generative adversarial networks. In *ICML*, 2017. **3**
- [16] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. **3**
- [17] K. Li, F. Xu, J. Wang, Q. Dai, and Y. Liu. A data-driven approach for facial expression synthesis in video. In *CVPR*, 2012. **2, 3**
- [18] M. Li, W. Zuo, and D. Zhang. Deep identity-aware transfer of facial attributes. *arXiv preprint arXiv:1610.05586*, 2016. **3**
- [19] X. Li, G. Mori, and H. Zhang. Expression-invariant face recognition with expression classification. In *CRV*, 2006. **1, 4**
- [20] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *CVPRW*, 2010. **5, 10**
- [21] I. Matthews and S. Baker. Active appearance models revisited. *IJCV*, 60(2):135–164, 2004. **1, 3, 4**
- [22] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. **3, 4**
- [23] U. Mohammed, S. J. D. Prince, and J. Kautz. Visio-ization: generating novel facial images. In *ACM SIGGRAPH*, 2009. **1, 3**
- [24] K. Olszewski, Z. Li, C. Yang, Y. Zhou, R. Yu, Z. Huang, S. Xiang, S. Saito, P. Kohli, and H. Li. Realistic dynamic facial textures from a single image using gans. In *ICCV*, 2017. **2**
- [25] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *BMVC*, 2015. **8**

- [26] F. Pighin, J. Hecker, D. Lischinski, R. Szeliski, and D. H. Salesin. Synthesizing realistic facial expressions from photographs. In *ACM SIGGRAPH Courses*, 2006. 1
- [27] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 3
- [28] S. Reed, K. Sohn, Y. Zhang, and H. Lee. Learning to disentangle factors of variation with manifold interaction. In *ICML*, 2014. 3
- [29] Z. Shu, E. Yumer, S. Hadap, K. Sunkavalli, E. Shechtman, and D. Samaras. Neural face editing with intrinsic image disentangling. In *CVPR*, 2017. 3
- [30] J. M. Susskind, G. E. Hinton, J. R. Movellan, and A. K. Anderson. Generating facial expressions with deep belief nets. In *Affective Computing*. 2008. 1, 3
- [31] B.-J. Theobald, I. Matthews, M. Mangini, J. R. Spies, T. R. Brick, J. F. Cohn, and S. M. Boker. Mapping and manipulating facial expression. *Language and speech*, 52(2-3):369–386, 2009. 1
- [32] J. Thies, M. Zollhofer, M. Stamminger, C. Theobald, and M. Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *CVPR*, 2016. 2
- [33] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *NIPS*, 2014. 3
- [34] X. Wu, R. He, and Z. Sun. A lightened cnn for deep face representation. *arXiv preprint arXiv:1511.02683*, 2015. 4, 8
- [35] F. Yang, L. Bourdev, E. Shechtman, J. Wang, and D. Metaxas. Facial expression editing in video using a temporally-smooth factorization. In *CVPR*, 2012. 1, 2, 3
- [36] F. Yang, J. Wang, E. Shechtman, L. Bourdev, and D. Metaxas. Expression flow for 3d-aware face component transfer. In *ACM Transactions on Graphics (TOG)*, 2011. 1, 2
- [37] R. Yeh, Z. Liu, D. B. Goldman, and A. Agarwala. Semantic facial expression editing using autoencoded flow. *arXiv preprint arXiv:1611.09961*, 2016. 1, 2
- [38] Z. Yi, H. Zhang, P. T. Gong, et al. Dualgan: Unsupervised dual learning for image-to-image translation. In *ICCV*, 2017. 3
- [39] Q. Zhang, Z. Liu, G. Quo, D. Terzopoulos, and H.-Y. Shum. Geometry-driven photorealistic facial expression synthesis. *TVCG*, 12(1):48–60, 2006. 1
- [40] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 1, 3

A. Supplementary Material

In this supplementary material, we present fully detailed information on 1) comparison experiment with ExprGAN [6] on the Oulu-CASIA dataset; 2) expression synthesis results for fine-grained control of eye status.

A.1. Comparison experiment with ExprGAN [6]

As mentioned in our paper, ExprGAN is the most similar work to ours. In this part, we take the same facial expression

synthesis experiments with ExprGAN [6] for comparison.

Fig. 8 shows the results of ExprGAN on the Oulu-CASIA dataset [4]. Images in the top row are input faces, and the rest are synthesized expressions with ascending intensities from top to down. Six types of expressions are synthesized, which correspond to the annotated expression classes in Oulu-CASIA dataset respectively. In addition, the neutral expression faces can also be generated in ExprGAN. In order to compare with ExprGAN, we take the same expression synthesis experiments and use the same input images with ExprGAN as shown in Fig. 9. Firstly, we employ the proposed expression removal model to recover the neutral expression faces, which are shown in the top row. Then, expression synthesis is conducted on neutral expression faces to generate various expressed faces. Following the setting in ExprGAN, we synthesize six types of expressions with five different intensities. It is worth noting that these two input images (subject id in the Oulu-CASIA dataset are P074 and P076 respectively) are not in our training dataset.

Due to the different way of image cropping, larger face areas (especially the chin areas) are covered in our experiments than ExprGAN. Particularly, the chin areas show wide variations along with expression changes, resulting in more difficulties in learning expression synthesis model. Comparing with ExprGAN, G2-GAN does much better in preserving identity information and keeping local details (such as hair in Fig. 9(a) and beard in Fig. 9(b)) through the expression transformation process. The neutral faces recovered by ExprGAN do not look like the ground truth images that are shown in the top row of Fig. 9, whereas G2-GAN is able to generate neutral faces without losing much identity information. Besides, we can see that the proposed G2-GAN can generate expressed faces with fine details such as wrinkles caused by frown and pout, whereas the results generated by ExprGAN tend to lack these details.

A.2. Expression synthesis with controlled eye status

The usage of face geometry in our framework provides an intuitive way for specifying target facial expression, with which we can generate special expressions such as “a lopsided grin with one eye open”. In this part, we show the ability of G2-GAN to synthesize facial expressions with controlled status of eye (open and closed). Since images in the Oulu-CASIA dataset are of low-resolution, we only take this experiment on the CK+ dataset [20].

Fig. 10 shows examples of synthesized images with various eye statuses. We can see that the proposed G2-GAN not only generates compelling perceptual results but also preserves identity information well. By directly manipulate the face geometry, we can perform fine-grained control of the eye status, which is hard for other generative model-based approaches. These results demonstrate the superiority of

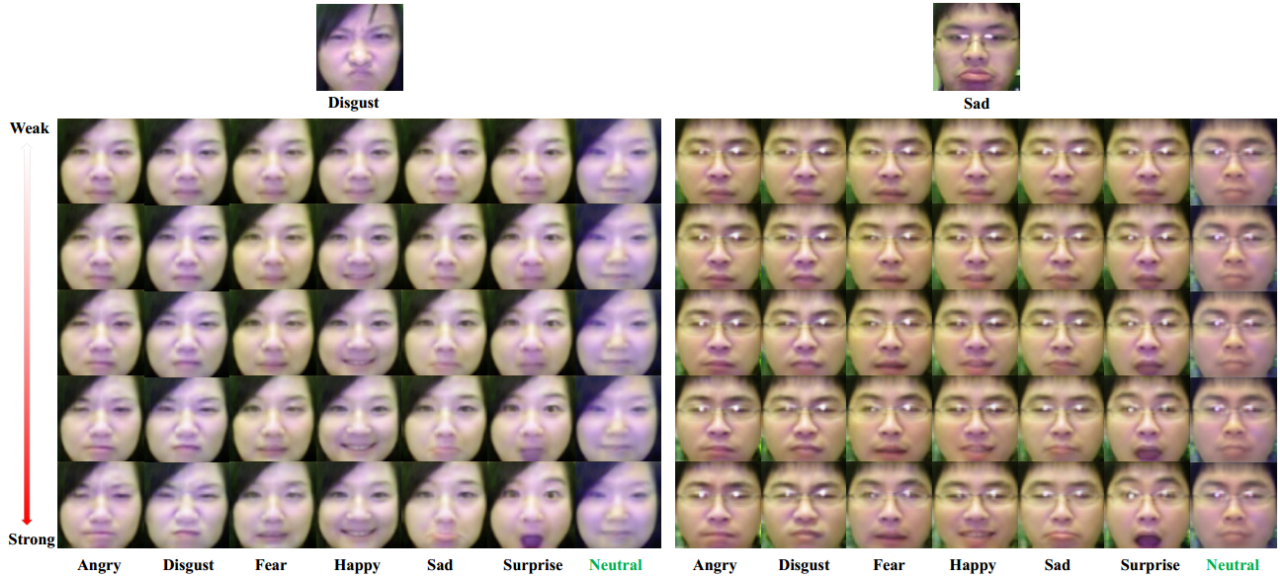


Figure 8. Facial expression synthesis results of ExprGAN [6] on the Oulu-CASIA dataset(originally shown in [6]).

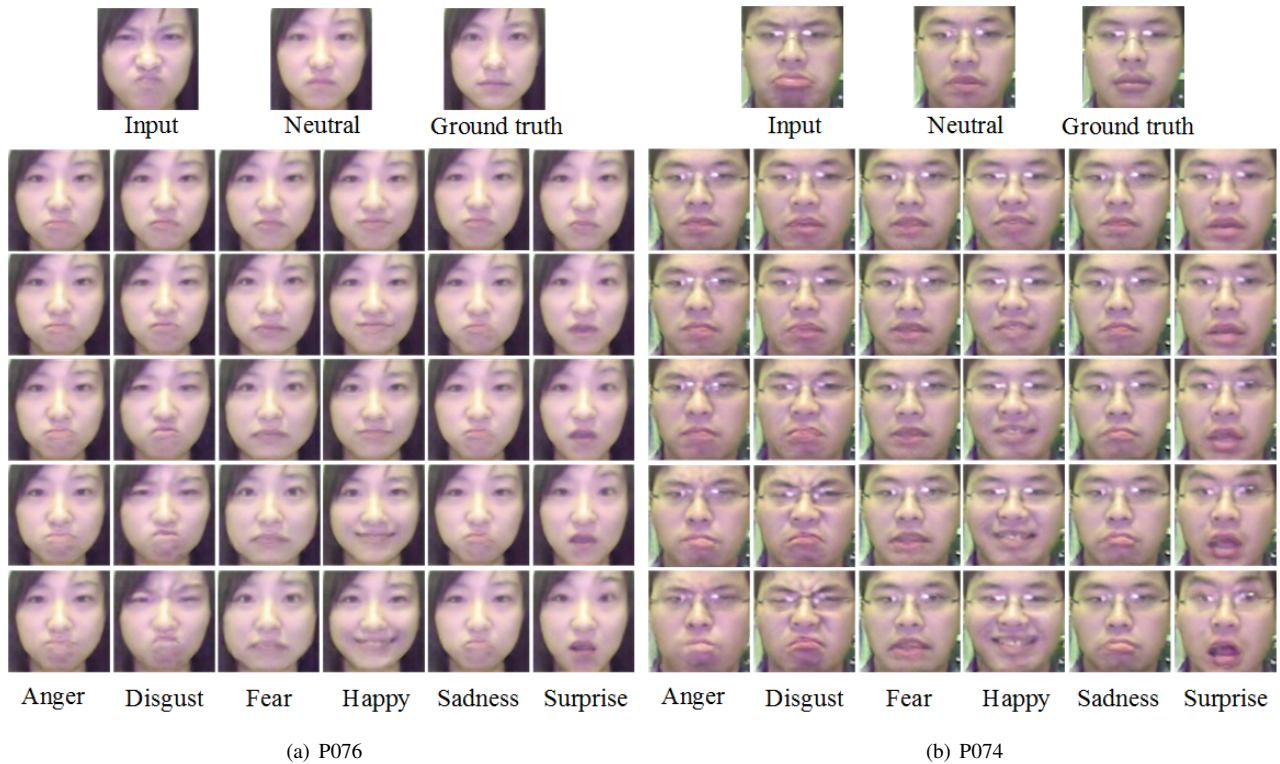


Figure 9. Facial expression synthesis results of the proposed G2-GAN. We take the same source images with ExprGAN [6] for comparison. The ‘neutral’ face image is generated from the ‘input’ face.

G2-GAN in operation-friendliness as well as the diversity of synthesized faces, suggesting potential applications for face edit. More results on the CK+ dataset and MultiPIE dataset [10] can be found in Fig. 11 and Fig. 12.



Figure 10. Examples of synthesized facial expressions with controlled status of eye. From left to right, eyes are gradually closed.



Figure 11. Examples of synthesized facial expressions on the CK+ dataset. Images in the first column are input faces, and the rest are input heatmaps and synthesized results.

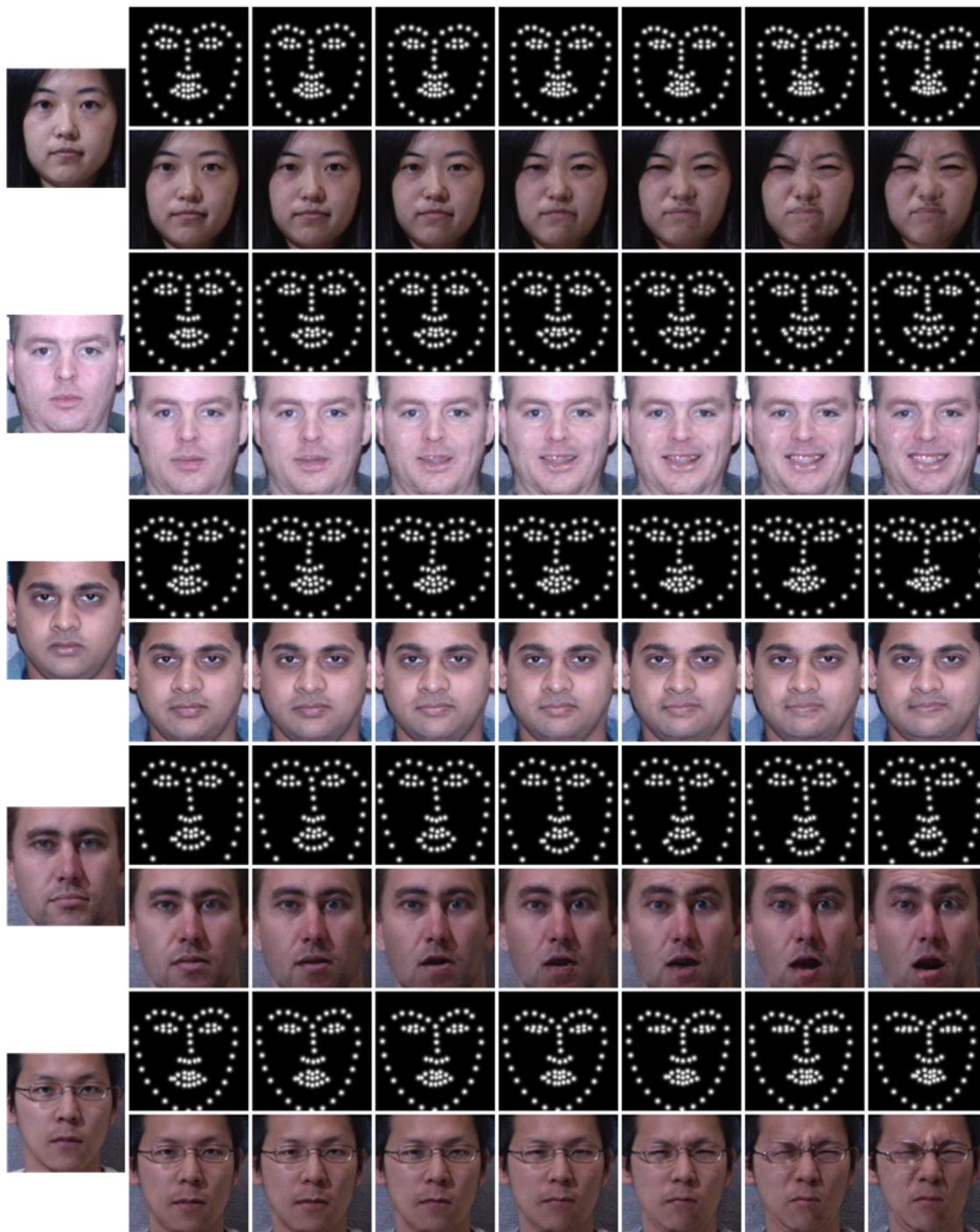


Figure 12. Examples of synthesized facial expressions on the MultiPIE dataset. Images in the first column are input faces, and the rest are input heatmaps and synthesized results.