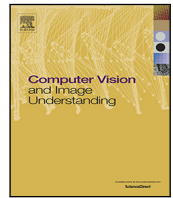




Contents lists available at ScienceDirect

Computer Vision and Image Understanding

journal homepage: www.elsevier.com/locate/cviu

Learn to synthesize and synthesize to learn[☆]

Behzad Bozorgtabar^{a,*}, Mohammad Saeed Rad^a, Hazim Kemal Ekenel^b, Jean-Philippe Thiran^{a,c}^a Signal Processing Laboratory (LTS5), École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland^b Istanbul Technical University, Istanbul, Turkey^c Department of Radiology, University Hospital Center (CHUV), University of Lausanne (UNIL), Lausanne, Switzerland

ARTICLE INFO

Keywords:

Attribute guided face image synthesis
Generative adversarial network
Facial expression recognition

ABSTRACT

Attribute guided face image synthesis aims to manipulate attributes on a face image. Most existing methods for image-to-image translation can either perform a fixed translation between any two image domains using a single attribute or require training data with the attributes of interest for each subject. Therefore, these methods could only train one specific model for each pair of image domains, which limits their ability in dealing with more than two domains. Another disadvantage of these methods is that they often suffer from the common problem of mode collapse that degrades the quality of the generated images. To overcome these shortcomings, we propose attribute guided face image generation method using a single model, which is capable to synthesize multiple photo-realistic face images conditioned on the attributes of interest. In addition, we adopt the proposed model to increase the realism of the simulated face images while preserving the face characteristics. Compared to existing models, synthetic face images generated by our method present a good photorealistic quality on several face datasets. Finally, we demonstrate that generated facial images can be used for synthetic data augmentation, and improve the performance of the classifier used for facial expression recognition.

1. Introduction

In this work, we are interested in the problem of synthesizing realistic faces by controlling the facial attributes of interest (e.g. expression, pose, lighting condition) without affecting the identity properties (see Fig. 1). In addition, this paper investigates learning from synthetic facial images for improving expression recognition accuracy. Synthesizing photo-realistic facial images has applications in human-computer interactions, facial animation and more importantly in facial identity or expression recognition. However, this task is challenging since image-to-image translation is ill-defined problem and it is difficult to collect images of varying attributes for each subject (e.g. images of different facial expressions for the same subject). The most notable solution is the incredible breakthroughs in generative models. In particular, Generative Adversarial Network (GAN) (Goodfellow et al., 2014) variants have achieved state-of-the-art results for the image-to-image translation task. These GAN models could be trained in both with paired training data (Isola et al., 2017) and unpaired training data (Kim et al., 2017; Zhu et al., 2017). Most existing GAN models (Shen and Liu, 2017; Zhu et al., 2017) are proposed to synthesize images of a single attribute, which make their training inefficient in the case of having multiple attributes, since for each attribute a separate model

is needed. In addition, GAN based approaches are often fragile in the common problem of mode collapse that degrades the quality of the generated images. To overcome these challenges, our objective is to use a single model to synthesize multiple photo-realistic images from the same input image with varying attributes simultaneously. Our proposed model, namely Learn to Synthesize and Synthesize to Learn (LSSL) is based on encoder–decoder structure, using the image latent representation, where we model the shared latent representation across image domains. Therefore, during inference step, by changing input face attributes, we can generate plausible face images owing attribute of interest. We introduce bidirectional learning for the latent representation, which we have found this loss term to prevent generator mode collapse. Moreover, we propose to use an additional face parsing loss to generate high-quality face images.

Our paper makes the following contributions:

1. This paper investigates domain adaptation using simulated face images for improving expression recognition accuracy. We show that how the proposed approach can be used to generate photo-realistic frontal facial images using synthetic face image and unlabeled real face images as the input. We compared our results

[☆] No author associated with this paper has disclosed any potential or pertinent conflicts which may be perceived to have impending conflict with this work. For full disclosure statements refer to <https://doi.org/10.1016/j.cviu.2019.04.010>.

* Corresponding author.

E-mail address: behzad.bozorgtabar@epfl.ch (B. Bozorgtabar).

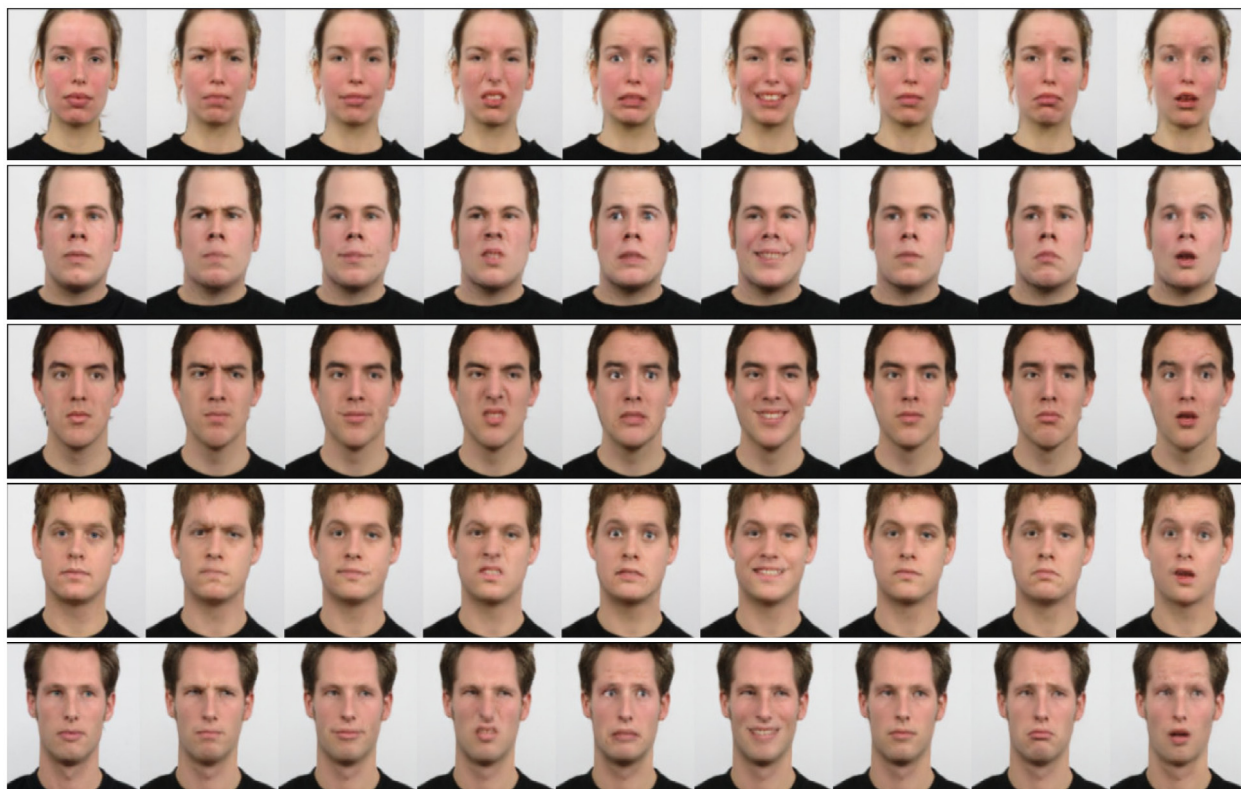


Fig. 1. Attribute guided facial image generation using LSSL on the Radboud Faces Database (RaFD) (Langner et al., 2010). The input neutral faces are fed into our model to exhibit specified attribute. Left to right: input neutral face and seven different attributes including angry, contemptuous, disgusted, fearful, happiness, neutral, sadness and surprised, respectively.

with SimGAN method (Shrivastava et al., 2017) in terms of expression recognition accuracy to see improvement in the realism of frontal faces. The source code is available at <https://github.com/CreativePapers/Learn-to-Synthesize-and-Synthesize-to-Learn>.

2. We show that use of our method leads to realistic generated images that contribute to improve the performance of expression recognition accuracy despite having small number of real training images. Further, compared to other variants of GAN models (Zhu et al., 2017; Perarnau et al., 2016; Choi et al., 2018), we show that a better performance can be attained through a proposed method to focus on the data augmentation process;
3. Unlike most of existing GAN based methods (Perarnau et al., 2016), which are trained with a large number of labeled and matching image pairs, the proposed method is adopted for unpaired image-to-image translation. As a matter of fact, the proposed method transfers the learnt characteristics between different classes;
4. The proposed method is capable of learning image-to-image translation among multiple domains using a single model. We introduce a bidirectional learning for the image latent representation to additionally enforce latent representation to capture shared features of different attribute categories and to prevent generator mode collapse. By doing so, we synthesize face photos with a desired attribute and translate an input image into another domain image.¹ Besides, we present face parsing loss and identity loss that help to preserve the face image local details and identity.

¹ We denote *domain* as a set of images owning the same attribute value.

2. Related work

Recently, GAN based models (Goodfellow et al., 2014) have achieved impressive results in many image synthesis applications, including image super-resolution (Ledig et al., 2017), image-to-image translation (pix2pix) (Isola et al., 2017) and CycleGAN (Zhu et al., 2017). We summarize contributions of few important related works in below:

Applications of GANs to face generation. Taigman et al. (2016) proposed a domain transfer network to tackle the problem of emoji generation for a given facial image. Lu et al. (2018) proposed attribute-guided face generation to translate low-resolution face images to high-resolution face images. Huang et al. (2017) proposed a Two-Pathway Generative Adversarial Network (TP-GAN) for photorealistic face synthesis by simultaneously considering local face details and global structures.

Image-to-image translation using GANs. Many of existing image-to-image translation methods e.g. Isola et al. (2017) and Shrivastava et al. (2017) formulated GANs in the supervised setting, where example image pairs are available. However, collecting paired training data can be difficult. On the other side, there are other GAN based methods, which do not require matching pairs of samples. For example, CycleGAN (Zhu et al., 2017) is capable to learn transformations from source to target domain without one-to-one mapping between two domain's training data. Li et al. (2016) proposed a Deep convolutional network model for Identity-Aware Transfer (DIAT) of the facial attributes. However, these GAN based methods could only train one specific model for each pair of image domains. Unlike the aforementioned approaches, we use a single model to learn to synthesize multiple photo-realistic images, each having specific attribute. More recently, IcGAN (Perarnau et al., 2016) and StarGAN (Choi et al., 2018) proposed image editing using AC-GAN (Odena et al., 2017) with conditional information. However, we use domain adaptation by adding the realism to the simulated faces

and there is no such a solution in these methods. Similar to Perarnau et al. (2016), Fader Networks Lample et al. (2017) proposed image synthesis model without needing to apply a GAN to the decoder output. However, these methods impose constraints on image latent space to enforce it to be independent from the attributes of interest, which may result in loss of information in generating attribute guided images.

GANS for facial frontalization and expression transfer. Zhang et al. (2018) proposed a method by disentangling the attributes (expression and pose) for simultaneous pose-invariant facial expression recognition and face images synthesis. Instead, we seek to learn attribute-invariant information in the latent space by imposing auxiliary classifier to classify the generated images. Qiao et al. (2018) proposed a Geometry-Contrastive Generative Adversarial Network (GC-GAN) for transferring continuous emotions across different subjects. However, this requires a training data with expression information, which may be expensive to obtain. Alternatively, our self-supervised approach automatically learns the required factors of variation by transferring the learnt characteristics between different emotion classes. Zhu et al. (2018) investigated GANs for data augmentation for the task of emotion classification. Lai and Lai (2018) proposed a multi-task GAN-based network that learns to synthesize the frontal face images from profile face images. However, they require paired training data of frontal and profile faces. Instead, we seek to add realism to the synthetic frontal face images without requiring real frontal face images during training. Our method could produce synthesis faces using synthetic frontal faces and real faces with arbitrary poses as input.

3. Methods

We first introduce our proposed multi-domain image-to-image translation model in Section 3.1. Then, we explain learning from simulated data by adding realism to simulated face images in Section 3.2. Finally, we discuss our implementation details and experimental results in Sections 4 and 5, respectively.

3.1. Learn to synthesize

Let \mathcal{X} and \mathcal{S} denote original image and side conditional image domains, respectively and \mathcal{Y} set of possible facial attributes, where we consider attributes including facial expression, head pose and lighting (see Fig. 2). As the training set, we have m triple inputs $(x_i \in \mathcal{X}, s_i \in \mathcal{S}, y_i \in \mathcal{Y})$, where x_i and y_i are the i th input face image and binary attribute, respectively and s_i represents the i th conditional side image as additional information to guide photo-realistic face synthesis. Then, for any categorical attribute vector y from the set of possible facial attributes \mathcal{Y} , the objective is to train a model that will generate photo-realistic version (x' or s') of the inputs (x and s) from image domains \mathcal{X} and \mathcal{S} with desired attributes y .

Our model is based on the encoder–decoder architecture with domain adversarial training. As the input to our expression synthesis method (see Fig. 3a), we propose to incorporate individual-specific facial shape model as the side conditional information s in addition to the original input image x . The shape model can be extracted from the configuration of the facial landmarks,² where the facial geometry varies with different individuals. Our goal is then to train a single generator G with encoder G_{enc} – decoder G_{dec} networks to translate the input pair (x, s) from source domains into their corresponding output images (x', s') in the target domain conditioned on the target domain attribute y and the inputs latent representation $G_{enc}(x, s)$, $G_{dec}(G_{enc}(x, s), y) \rightarrow x', s'$. The encoder $G_{enc} : (\mathcal{X}^{source}, \mathcal{S}^{source}) \rightarrow \mathbb{R}^{n \times \frac{h}{16} \times \frac{w}{16}}$ is a fully convolutional neural network with parameters θ_{enc} that encodes the input images into a low-dimensional feature space $G_{enc}(x, s)$, where n, h, w are the number of the feature channels and the input images dimensions,

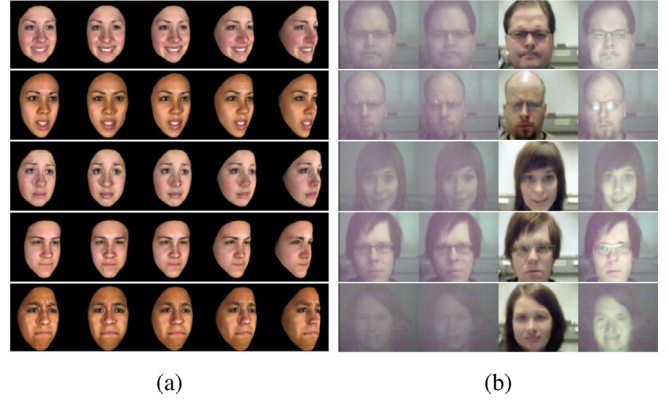


Fig. 2. Examples of facial attribute transfer. (a) Generating images with varying poses ranging from 0 to 45 degrees (yaw angle) in 15 degrees steps. (b) Generating face images with three different lighting conditions using face image with normal illumination as input: normal illumination (reconstruction), weak illumination and dark illumination, respectively.

respectively. The decoder $G_{dec} : (\mathbb{R}^{n \times \frac{h}{16} \times \frac{w}{16}}, \mathcal{Y}) \rightarrow \mathcal{X}^{target}, \mathcal{S}^{target}$ is the sub-pixel (Shi et al., 2016) convolutional neural network with parameters θ_{dec} that produce realistic images with target domain attribute y and given the latent representation $G_{enc}(x, s)$. The precise architectures of the neural networks are described in Section 4.1. During training, we randomly use a set of target domain attributes y to make the generator more flexible in synthesizing images. In the following, we introduce the objectives for the proposed model optimization.

GAN loss. We introduce a model that discovers cross-domain image translation with GANs. Moreover, at the inference time, we should be able to generate diverse facial images by only changing attribute of interest. By doing so, we seek to learn attribute-invariant information in the latent space representing the shared features of the images sampled for different attributes. It means if the original and target domains are semantically similar (e.g. facial images of different expressions), we expect the common features across domains to be captured by the same latent representation. Then, the decoder must use the target attribute to perform image-to-image translation from the original domain to the target domain. However, this learning process is unsupervised as for each training image from the source domain, its counterpart image in the target domain with attribute y is unknown. Therefore, we propose to train an additional neural network called the discriminator D (with the parameters θ_{dis}) using an adversarial formulation to not only distinguish between real and fake generated images, but also to classify the image to its corresponding attribute categories. We use Wasserstein GAN (Gulrajani et al., 2017) objective with a gradient penalty loss \mathcal{L}_{gp} (Arjovsky et al., 2017) formulated as below:

$$\mathcal{L}_{GAN} = \mathbb{E}_{x,s} [D_{src}(x, s)] - \mathbb{E}_{x,s,y} [D_{src}(G_{dec}(G_{enc}(x, s), y))] - \lambda_{gp} \mathcal{L}_{gp}(D_{src}), \quad (1)$$

The term $D_{src}(\cdot)$ denotes a probability distribution over image sources given by D . The hyper-parameter λ_{gp} is used to balance the GAN objective with the gradient penalty. A generator (encoder–decoder networks) used in our model has to play two roles: learning the attribute invariance representation for the input images and is trained to maximally fool the discriminator in a *min–max* game. On the other hand, the discriminator simultaneously seeks to identify the fake examples for each attribute.

Attribute classification loss. We deploy a classifier by returning additional output from the discriminator to perform an auxiliary task of classifying the synthesized and real facial images into their respective

² We use dlib regression trees algorithm.

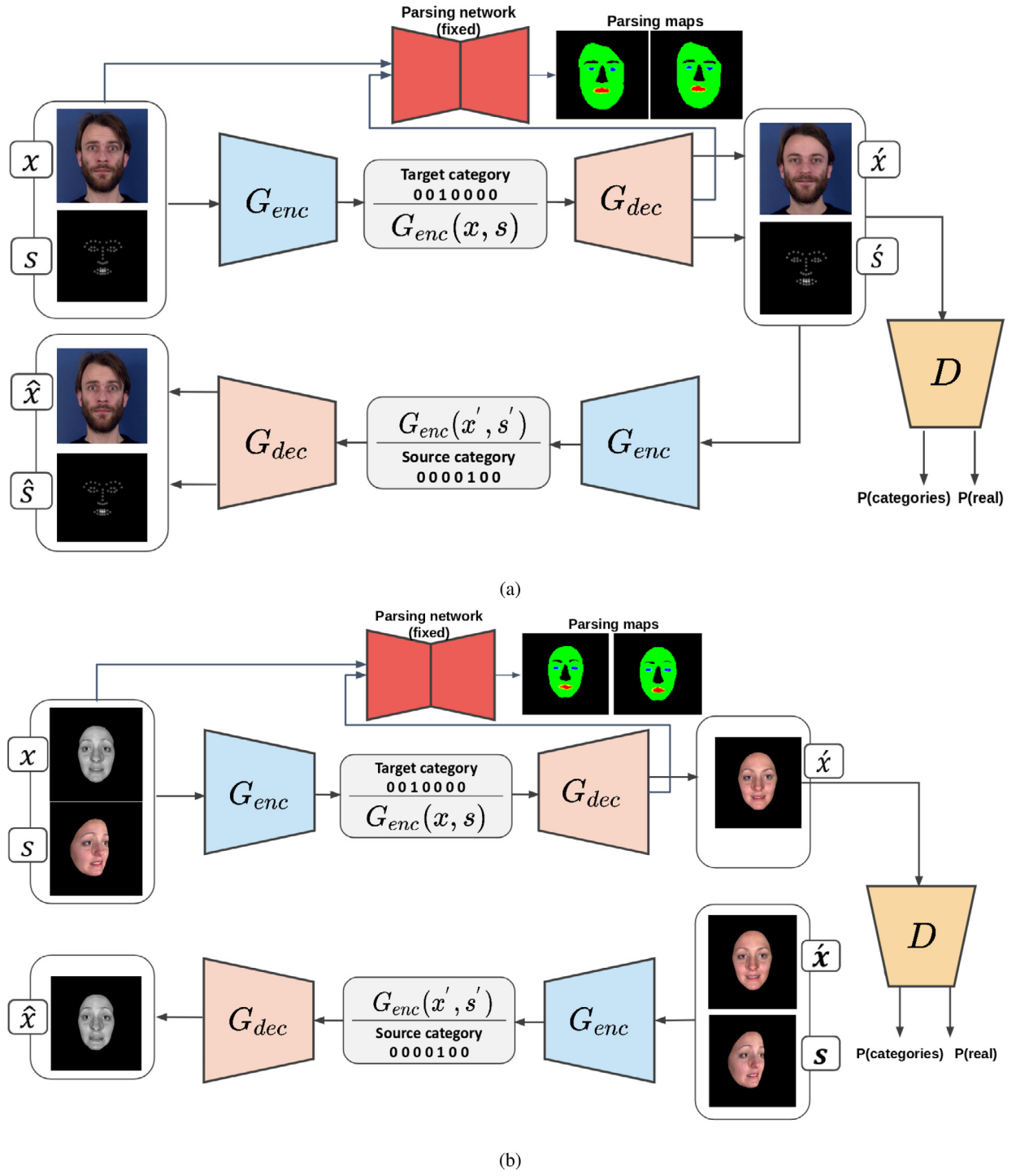


Fig. 3. Overview of our proposed LSSL method. (a) Attribute-guided face image synthesis model, consisting of three networks, an encoder–decoder generator G , a discriminator D and face parsing network P . A discriminator's job is to discriminate the realism of synthetic pair images and to guarantee correct attribute classification on the generated face images. (b) Pose normalization model, which takes a synthetic face image and unlabeled real image as input and generates photo-realistic version of the frontal face through identical networks used in (a). The only difference is that a discriminator takes only one single image as input.

attribute categories. An attribute classification loss of real images \mathcal{L}_{cls_r} to optimize the discriminator parameters θ_{dis} is defined as follow:

$$\min_{\theta_{dis}} \mathcal{L}_{cls_r} = \mathbb{E}_{x,s,y'} [\ell_r(x, s, y')],$$

$$\ell_r(x, s, y') = \sum_{i=1}^m -y_i' \log D_{cls}(x, s) - (1 - y_i') \log (1 - D_{cls}(x, s)),$$
(2)

Here, y' denotes original attributes categories for the real images. ℓ_r is the summation of binary cross-entropy losses of all attributes. Besides,

an attribute classification loss of fake images \mathcal{L}_{cls_f} used to optimize the generator parameters ($\theta_{enc}, \theta_{dec}$), formulated as follow:

$$\min_{\theta_{enc}, \theta_{dec}} \mathcal{L}_{cls_f} = \mathbb{E}_{x',s',y} [\ell_f(x', s', y)],$$

$$\ell_f(x', s', y) = \sum_{i=1}^m -y_i \log D_{cls}(x', s') - (1 - y_i) \log (1 - D_{cls}(x', s')),$$
(3)

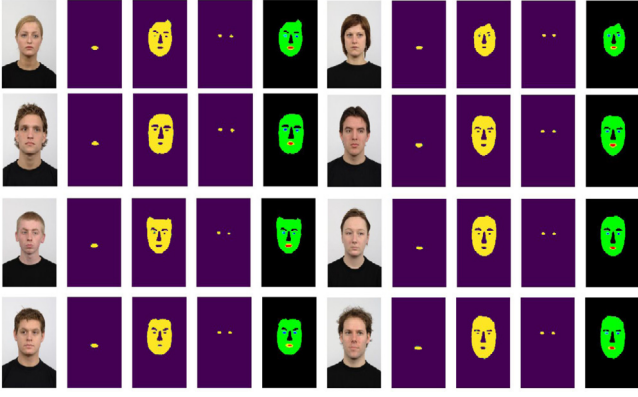


Fig. 4. Face parsing maps on the RaFD dataset. **Left to right:** input *neutral* face and parsing maps for its constituent facial parts containing lips (second column), face skin (third column), eyes (fourth column) and color visualization generated by all three category parsing maps (last column), respectively. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

where x' and s' are the generated images and auxiliary outputs, which should correctly own the target domain attributes y . \mathcal{L}_f denotes summing up the cross-entropy losses of all fake images.

Identity loss. Using the identity loss, we aim to preserve the attribute-excluding facial image details such as facial identity before and after image translation. By doing so, we use a pixel-wise l_1 loss to enforce the details consistency of the face original domain and suppress the face blurriness:

$$\mathcal{L}_{id} = \mathbb{E}_{x,s,y'} \left[\left\| G_{dec} (G_{enc} (x, s), y) - x \right\|_1 \right], \quad (4)$$

Face parsing loss. The face important components (e.g., lips and eyes) are typically small and cannot be well reconstructed by solely minimizing the identity loss on the whole face image. Therefore, we use



Fig. 6. Facial attribute transfer results of LSSL compared with IcGAN (Perarnau et al., 2016) and CycleGAN (Zhu et al., 2017), respectively.

a face parsing loss to further improve the harmony of the synthetic faces. As our face parsing network, we use U-Net (Ronneberger et al., 2015) trained on the Helen dataset (Le et al., 2012), which has ground truth face semantic labels, for training parsing network. Instead of utilizing all semantic labels, we use three key face components (lips, eyes and face skin). Once the network is trained, it remains fixed in our framework. The parsing loss is back-propagated to the generator to further regularize generator. Fig. 4 shows some parsing results on the RaFD dataset (Langner et al., 2010).

$$\mathcal{L}_p = \mathbb{E}_{x,s,y'} [A_p (P(x) - P(x'))], \quad (5)$$

where $A_p(\cdot, \cdot)$ denotes a function to compute pixel-wise softmax loss and $P(\cdot)$ is the face parsing network.

Bidirectional loss. Using GAN loss alone usually leads to mode collapse, generating identical labels regardless of the input face photo.

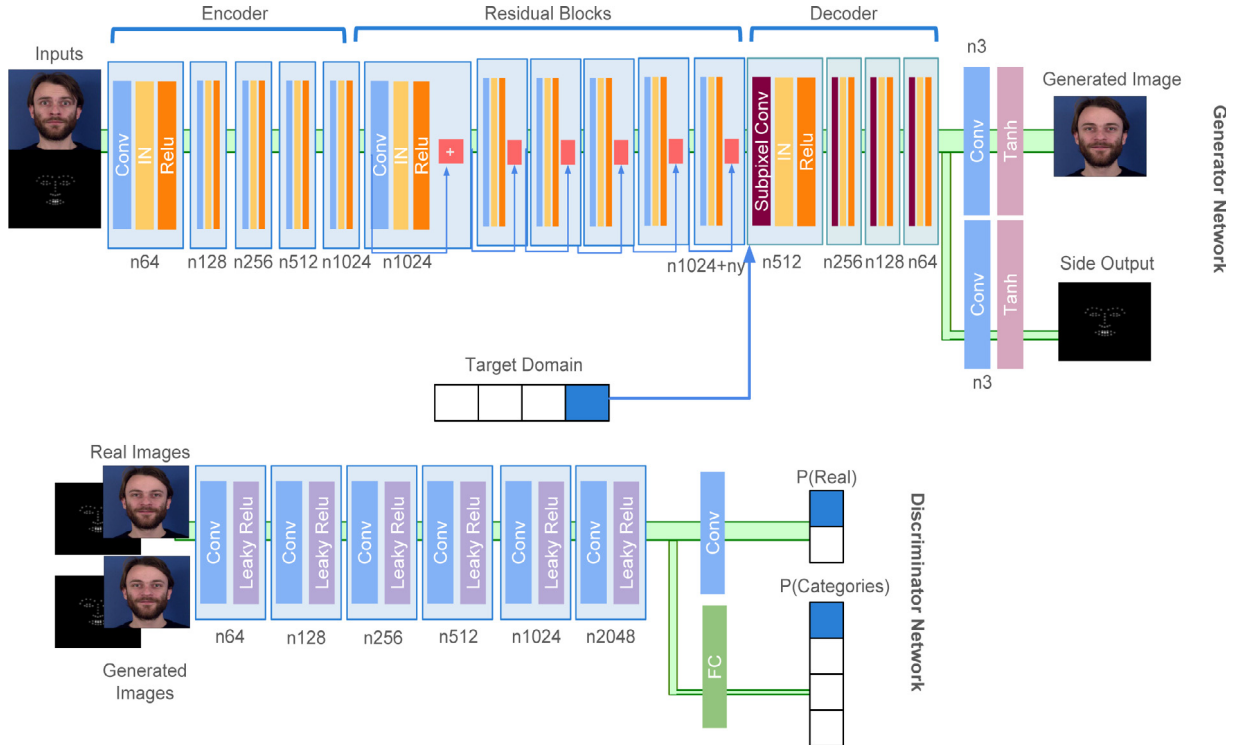


Fig. 5. Architecture of generator (top) and discriminator network (bottom). n_y denotes the dimension of domain attributes. IN and n and FC denote instance normalization, number of feature maps and fully connected layer, respectively.

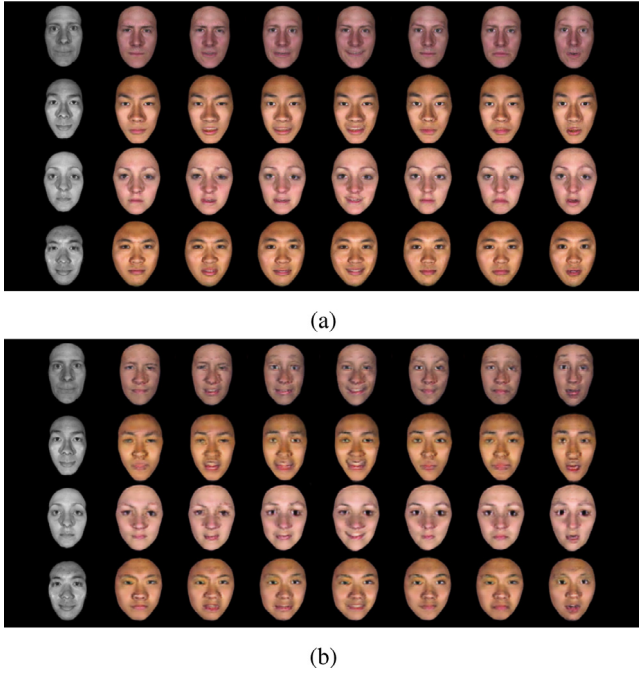


Fig. 7. Pose-normalized face attribute transfer results of (a) LSSL method compared with (b) SimGAN method (Shrivastava et al., 2017) on the BU-3DFE dataset (Yin et al., 2006). The input synthetic frontal face and real profile face are fed into our model to exhibit specified attribute. Left to right: input synthetic face and seven different attributes including *angry*, *disgusted*, *fearful*, *happiness*, *neutral*, *sadness* and *surprised*, respectively.

This problem has been observed in various applications of conditional GANs (Isola et al., 2017; Dosovitskiy and Brox, 2016) and to our knowledge, there is still no proper approach to deal with this issue. To address this problem, we show that using the trained generator, images of different domains can be translated bidirectionally. We decompose this objective into two terms: a bidirectional loss for the image latent representation, and a bidirectional loss between synthesized images and original input images, respectively. This objective is formulated using l_1 loss as follow:

$$\begin{aligned} \mathcal{L}_{bi} = & \mathbb{E}_{x,s,y'} [\|x - \hat{x}\|_1 + \|s - \hat{s}\|_1] + \\ & \mathbb{E}_{x,s,y} [\|G_{enc}(x,s) - G_{enc}(x',s')\|_1], \end{aligned} \quad (6)$$

$$\begin{aligned} x', s' = & G_{dec}(G_{enc}(x,s), y), \\ \hat{x}, \hat{s} = & G_{dec}(G_{enc}(x',s'), y'), \end{aligned}$$

In the above equation, \hat{x} and \hat{s} denote the reconstructed original image and the side conditional image, respectively. Unlike (Zhu et al., 2017), where only the cycle consistency losses are used at the image level, we additionally seek to minimize the reconstruction loss using latent representation.

Overall objective. Finally, the generator G is trained with a linear combination of five loss terms: adversarial loss, attribute classification loss for the fake images, bidirectional loss, identity loss and face parsing loss. Meanwhile, the discriminator D is optimized using an adversarial loss and attribute classification loss for the real images:

$$\begin{aligned} \mathcal{L}_G = & \mathcal{L}_{GAN} + \lambda_{bi}\mathcal{L}_{bi} + \lambda_{cls}\mathcal{L}_{cls_f} + \lambda_{id}\mathcal{L}_{id} + \lambda_p\mathcal{L}_p, \\ \mathcal{L}_D = & -\mathcal{L}_{GAN} + \lambda_{cls}\mathcal{L}_{cls_r}, \end{aligned} \quad (7)$$

where λ_{bi} , λ_p , λ_{id} and λ_{cls} are hyper-parameters, which tune the importance of bidirectional loss, face parsing loss, identity loss and attribute classification loss, respectively.

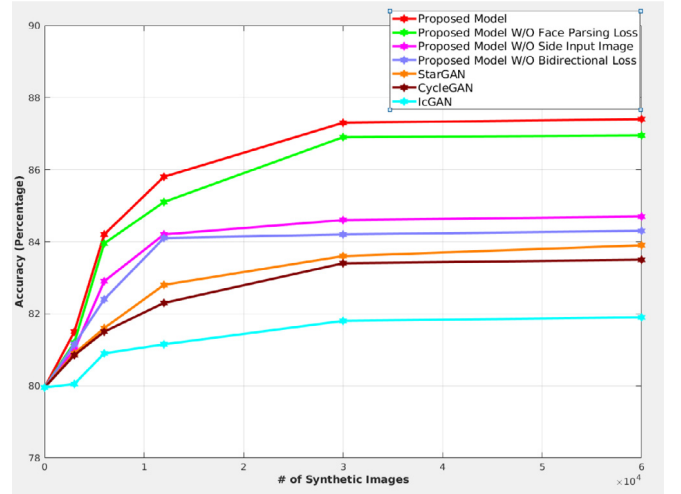


Fig. 8. Impact of the amount of training synthetic images on performance in terms of expression recognition accuracy.

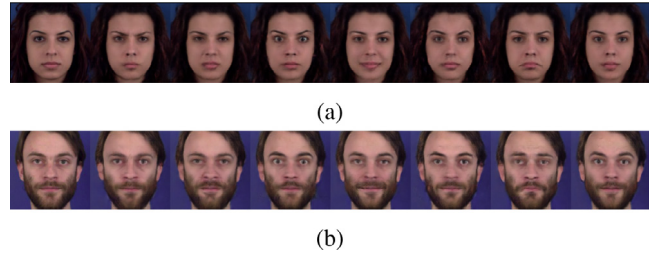


Fig. 9. Facial attribute transfer results from our proposed method for (a) subject 1 and (b) subject 20, respectively. The input face images are manipulated to exhibit desired attribute. Left to right: input *neutral* face and seven different attributes including *anger*, *disgust*, *fear*, *happiness*, *neutral*, *sadness* and *surprise*, respectively.

3.2. Synthesize to learn

In an unconstrained face expression recognition, accuracy will drop significantly for large pose variations. The key solution would be using simulated faces rendered in frontal view. However, learning from synthetic face images can be problematic due to a distribution discrepancy between real and synthetic images. Here, our proposed model generates realistic face images given real profile face with arbitrary pose and a simulated face image as input (see Fig. 3b). We utilize a 3D Morphable Model using bilinear face model (Vlasic et al., 2005) to construct a simulated frontal face image from multiple camera views. Here, the discriminator’s role is to discriminate the realism of synthetic face images using unlabeled real profile face images as a conditional side information. In addition, using the same discriminator, we can generate face images exhibiting different expressions.

We compare the results of LSSL with SimGAN method (Shrivastava et al., 2017) on the BU-3DFE dataset (Yin et al., 2006) to evaluate the realism of face images. SimGAN method Shrivastava et al. (2017) considers learning from simulated and unsupervised images through adversarial training. However, SimGAN is devised for much simpler scenarios e.g., eye image refinement. In addition, categorical information was ignored in SimGAN, which limits the model generalization. In contrast, LSSL overcomes this issue by introducing attribute classification loss into objective function. For a fair comparison with SimGAN method, we add the attribute classification loss by modifying the SimGAN’s discriminator, while keeping the rest of network unchanged. We achieve more visually pleasing results on test data compared to the SimGAN method (see Fig. 7).

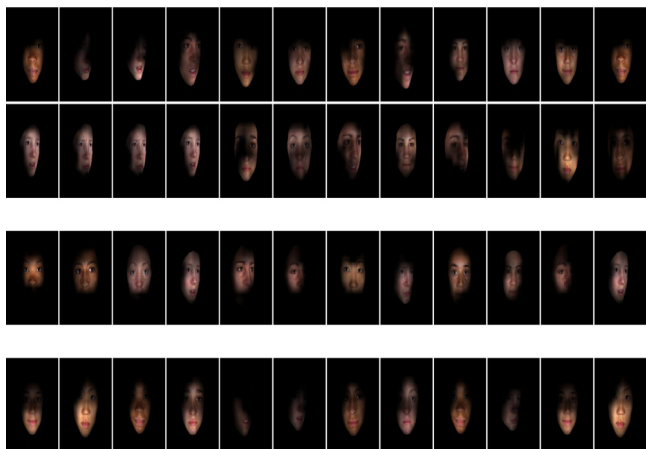


Fig. 10. Visualization of some hidden units in the encoder of LSSL trained on the BU-3DFE dataset (Yin et al., 2006). We highlight regions of face images that a particular convolutional hidden unit maximally activates on.

4. Implementation details

All networks are trained using Adam optimizer (Kingma and Ba, 2014) ($\beta_1 = 0.5, \beta_2 = 0.999$) and with a base learning rate of 0.0001. We linearly decay learning rate after the first 100 epochs. We use a simple data augmentation with only flipping the images horizontally. The input image size and the batch size are set to 128×128 and 8 for all experiments, respectively. We update the discriminator five times for each generator (encoder–decoder) update. The hyper-parameters in Eqs. (7) and (1) are set as: $\lambda_{bj} = 10$ and $\lambda_{id} = 10$, $\lambda_p = 10$, $\lambda_{gp} = 10$ and $\lambda_{cls} = 1$, respectively. The whole model is implemented using PyTorch on a single NVIDIA GeForce GTX 1080.

4.1. Networks architectures

For the discriminator, we use PatchGAN (Isola et al., 2017) that penalizes structure at the scale of image patches. In addition, LSSI has the generator network composed of five convolutional layers with the stride size of two for downsampling, six residual blocks, and four transposed convolutional layers with the stride size of two for upsampling. We use sub-pixel convolution instead of transposed convolution followed by instance normalization (Ba et al., 2016). For the face parsing network, we used the same net architecture as U-Net proposed in Ronneberger et al. (2015), but our face parsing network consists of depthwise convolutional blocks proposed by MobileNets (Sandler et al., 2018). The network architecture of LSSI is shown in Fig. 5.

5. Experimental results

In this section, we first propose to carry out comparison between our LSSI method and recent methods in image-to-image translation from a qualitative perspective, then we demonstrate the generality of our method (quantitative analysis) using different techniques for the face expression recognition.

5.1. Datasets

Oulu-CASIA VIS (Zhao et al., 2011): This dataset contains 480 sequences (from 80 subjects) of six basic facial expressions under the visible (VIS) normal illumination conditions. The sequences start from a neutral face and end with peak facial expression. This dataset is chosen due to high intra-class variations caused by the personal attributes. We conducted our experiments using subject-independent 10-fold cross-validation strategy.

MUG (Aifanti et al., 2010): The MUG dataset contains image sequences of seven different facial expressions belonging to 86 subjects comprising 51 men and 35 women. The image sequences were captured with a resolution of 896×896 . We used image sequences of 52 subjects and the corresponding annotation, which are available publicly via the internet.

BU-3DFE (Yin et al., 2006): The Binghamton University 3D Facial Expression Database (BU-3DFE) (Yin et al., 2006) contains 3D models from 100 subjects, 56 females and 44 males. The subjects show a neutral face as well as six basic facial expressions and at four different intensity levels. Following the setting in Tariq et al. (2013) and Zhang et al. (2018), we used an OpenGL based tool from the database creators to render multiple views from 3D models in seven pan angles ($0^\circ, \pm 15^\circ, \pm 30^\circ, \pm 45^\circ$).

RaFD (Langner et al., 2010): The Radboud Faces Database (RaFD) contains 4,824 images belonging to 67 participants. Each subject makes eight facial expressions.

Qualitative evaluation. As shown in Fig. 6, our facial attribute transfer test results (unseen images during the training step) are more visually pleasing compared to recent baselines including IcGAN (Perarnau et al., 2016) and CycleGAN. (Zhu et al., 2017). We believe that our proposed losses (parsing loss and identity losses) help to preserve the face image details and identity. IcGAN even fails to generate subjects with desired attributes, while our proposed method could learn attribute invariant features applicable to synthesize multiple images with desired attributes. In addition, to evaluate the proposed pose normalization method, the face attribute transfer results of our proposed method have been compared with the SimGAN method (Shrivastava et al., 2017) on the BU-3DFE dataset (Yin et al., 2006) (see Fig. 7).

Quantitative evaluation. To conduct the quantitative analysis, we apply LSSI to data augmentation for facial expression recognition. We augment real images from Oulu-CASIA VIS dataset with the synthetic expression images generated by LSSI as well as its variants and then compare with other methods to train an expression classifier. The purpose of this experiment is to introduce more variability and enrich the dataset further, in order to improve the expression recognition performance. In particular, from each of the six expression category, we generate 0.5k, 1k, 2k, 5k and 10k images, respectively. As shown in Fig. 8, when the number of synthetic images is increased to 30k, the accuracy is improved drastically, reaching to 87.40%. The performance starts to become saturated when more images (60k) are used. We achieved a higher recognition accuracy value using the images generated from LSSI than other CNN-based methods including popular generative model, StarGAN (Choi et al., 2018) (see Table 1). This suggests that our model has learned to generate more realistic facial images controlled by the expression category. In addition, we evaluate the sensitivity of the results for different components of LSSI method (face parsing loss, bidirectional loss and side conditional image, respectively). We observe that our LSSI method trained with each of the proposed loss terms yields a notable performance gain in facial expression recognition.

Moreover, we evaluate the performance of LSSI on the MUG facial expression dataset (Aifanti et al., 2010) using the video frames of the peak expressions. Fig. 9 shows sample facial attribute transfer results on the MUG facial dataset (Aifanti et al., 2010). It should be noted that the MUG facial expression dataset are only available to authorized users. We only have permission from few subjects including 1 and 20 for using their photos in our paper. In Table 2, we report the results of average accuracy of a facial expression on synthesized images. We trained a facial expression classifier with (90%/10%) splitting for training and test sets using a ResNet-50 (He et al., 2016), resulting in a near-perfect accuracy of 90.42%. We then trained each of baseline models including CycleGAN, IcGAN and StarGAN using the same training set and performed image-to-image translation on the same test set. Finally,

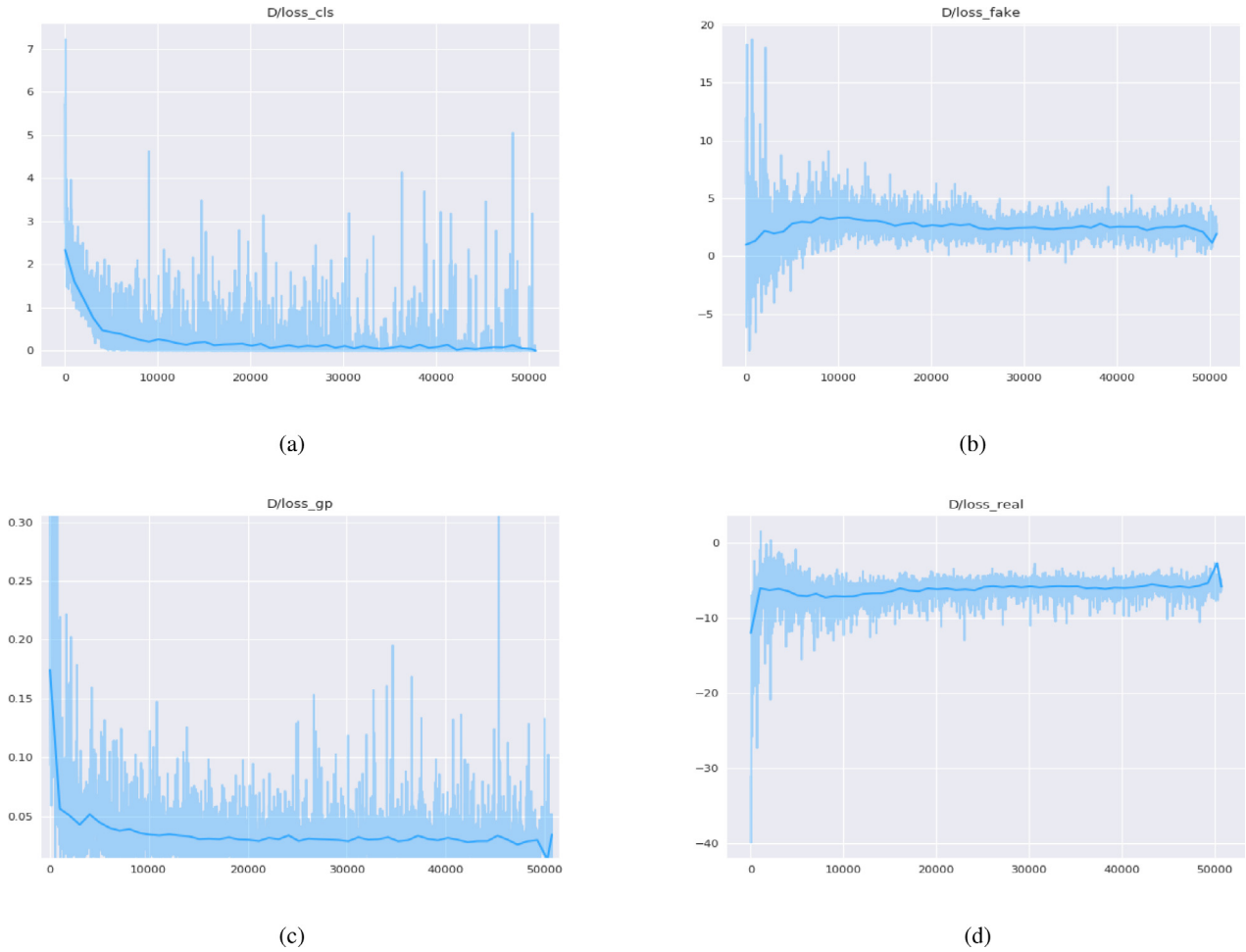


Fig. 11. The training losses for the network’s discriminator on the RaFD dataset (Langner et al., 2010). (a) An attribute classification loss of real images, (b) the discriminator loss for the pair of fake generated images, (c) discriminator gradient penalty loss and (d) discriminator loss for the pair of real images, respectively.

Table 1

Performance comparison of expression recognition accuracy between the proposed method and other state-of-the-art methods.

Method	Accuracy
HOG 3D (Klaser et al., 2008)	70.63%
AdaLBP (Zhao et al., 2011)	73.54%
Atlases (Guo et al., 2012)	75.52%
STM-ExpLet (Liu et al., 2014)	74.59%
DTAGN (Jung et al., 2015)	81.46%
StarGAN (Choi et al., 2018)	83.90%
LSSL W/O side input	84.70%
LSSL W/O bidirectional loss	84.30%
LSSL W/O face parsing loss	86.95%
LSSL	87.40%

we classified the expression of these generated images using the above-mentioned classifier. As can be seen in Table 2, our model achieves the highest classification accuracy (close to real image), demonstrating that our model could generate the most realistic expressions among all the methods compared.

Pose normalization analysis. Using BU-3DFE dataset (Yin et al., 2006), we have designed subject-independent experimental setup. We performed 5-fold cross validation using 100 subjects. Training data includes images of 80 (frontal face) subjects, while test data includes images of 20 subjects with varying poses. We use VGG-Face model (Parkhi et al., 2015), which is pretrained on the (RaFD) (Langner et al., 2010) and then we further fine-tune it on the frontal face images from

Table 2

Performance comparison on the MUG dataset in terms of average classification accuracy.

Method	Accuracy
Real test set	90.42%
CycleGAN (Zhu et al., 2017)	84.40%
IcGAN (Perarnau et al., 2016)	80.32%
StarGAN (Choi et al., 2018)	85.15%
LSSL W/O face parsing loss	89.91%
LSSL	90.35%

BU-3DFE dataset. It can be observed from Table 3 that pose normalization helps to improve expression recognition performance of the non-frontal faces (ranging from 15 to 45 degrees in 15 degrees steps). Having said that, adding realism to simulated face images helps to bring additional gains in terms of expression recognition accuracy. In particular, our method outperforms two recent works, Lai and Lai (2018) and Zhang et al. (2018) that addressed pose normalization task. Our proposed losses (parsing loss and identity losses) facilitates the synthesized frontal face images to preserve much detail of face characteristics (e.g. expression and identity).

5.2. Visualizing representation

Fig. 10 visualizes some activations of hidden units in the fifth layer of an encoder (the first component of the generator). Although all units are not semantic, but these visualizations indicate that the

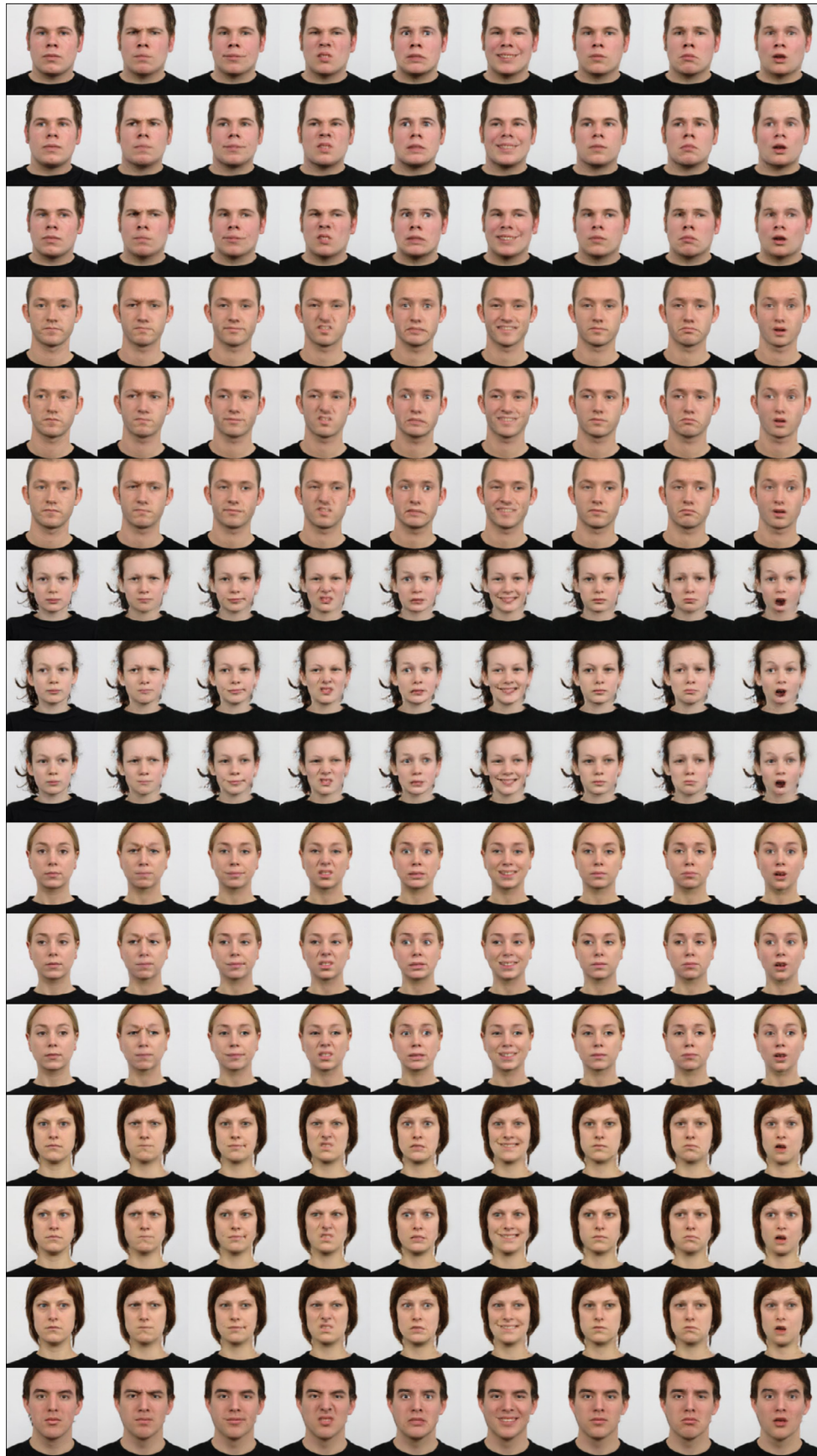


Fig. 12. Facial expression synthesis on the Radboud Faces Database (RaFD). Left to right: input *neutral* face and synthesis results of all eight emotion classes including *angry*, *contemptuous*, *disgusted*, *fearful*, *happiness*, *neutral*, *sadness* and *surprised*, respectively.

Table 3
Recognition accuracies on normalized face images at different pose angles.

Method	± 15	± 30	± 45
Real profile face images	70.15%	66.50%	58.90%
Simulated frontal face images	70.91%	65.90%	59.30%
CycleGAN	71.60%	67.32%	61.50%
Lai and Lai (2018)	71.45%	67.60%	61.95%
Zhang et al. (2018)	71.72%	67.65%	62.10%
LSSL W/O face parsing loss	72.10%	68.52%	63.35%
LSSL	72.70%	69.10%	64.05%

network learns to identify the most informative visual cues from the face regions.

5.3. Training losses additional qualitative results

Fig. 11 shows the training losses of the proposed attribute guided face image synthesis model for the discriminator. Here, we use the face landmark heatmap as the side conditional image. The face landmark heatmap contains 2D Gaussians centered at the landmarks' locations, which are then concatenated with the input image to synthesize different facial expressions on the RaFD dataset (Langner et al., 2010). In addition, the target attribute label is spatially replicated and concatenated with the latent feature. Results in Fig. 11 are for 100 epochs, 50,000 iterations of training on the RaFD dataset. Moreover, Fig. 12 shows additional images generated by LSSL.

6. Conclusion

In this work, we introduced LSSL, a model for multi-domain image-to-image translation applied to the task of face image synthesis. We present attribute guided face image generation to transform a given image to various target domains controlled by desired attributes. We argue that learning image-to-image translation between image domains requires a proper modeling of the shared latent representation across image domains. Additionally, we proposed face parsing loss and identity loss to preserve much detail of face characteristics (e.g. identity). More importantly, we seek to add realism to the synthetic images while preserving the face pose angle. We also demonstrate that the synthetic images generated by our method can be used for data augmentation to enhance facial expression classifier's performance. We reported promising results on the task of domain adaptation by adding the realism to the simulated faces. We showed that by leveraging the synthetic face images as a form of data augmentation, we can achieve significantly higher average accuracy compared with the state-of-the-art result.

Acknowledgment

This work is supported by the European Union's Horizon 2020 research and innovation programme under grant agreement No 688900 (ADAS&ME project - <http://www.adasandme.com>).

References

Aifanti, N., Papachristou, C., Delopoulos, A., 2010. The MUG facial expression database. In: *Image Analysis for Multimedia Interactive Services (WIAMIS)*, 2010 11th International Workshop on. IEEE, pp. 1–4.

Arjovsky, M., Chintala, S., Bottou, L., 2017. Wasserstein generative adversarial networks. In: *International Conference on Machine Learning*. pp. 214–223.

Ba, J.L., Kiros, J.R., Hinton, G.E., 2016. Layer normalization, arXiv preprint arXiv:1607.06450.

Choi, Y., Choi, M., Kim, M., Ha, J.-W., Kim, S., Choo, J., 2018. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 8789–8797.

Dosovitskiy, A., Brox, T., 2016. Generating images with perceptual similarity metrics based on deep networks. In: *Advances in Neural Information Processing Systems*. pp. 658–666.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets. In: *Advances in Neural Information Processing Systems*. pp. 2672–2680.

Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C., 2017. Improved training of wasserstein gans. In: *Advances in Neural Information Processing Systems*. pp. 5769–5779.

Guo, Y., Zhao, G., Pietikäinen, M., 2012. Dynamic facial expression recognition using longitudinal facial expression atlases. In: *Computer Vision—ECCV 2012*. Springer, pp. 631–644.

He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 770–778.

Huang, R., Zhang, S., Li, T., He, R., et al.,

Isola, P., Zhu, J.-Y., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1125–1134.

Jung, H., Lee, S., Park, S., Lee, I., Ahn, C., Kim, J., 2015. Deep temporal appearance-geometry network for facial expression recognition. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV 2015)*. IEEE.

Kim, T., Cha, M., Kim, H., Lee, J.K., Kim, J., 2017. Learning to discover cross-domain relations with generative adversarial networks. In: *International Conference on Machine Learning*. pp. 1857–1865.

Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980.

Klaser, A., Marszałek, M., Schmid, C., 2008. A spatio-temporal descriptor based on 3d-gradients. In: *BMVC 2008-19th British Machine Vision Conference*. British Machine Vision Association, pp. 275–1.

Lai, Y.-H., Lai, S.-H., 2018. Emotion-preserving representation learning via generative adversarial network for multi-view facial expression recognition. In: *Automatic Face & Gesture Recognition (FG 2018)*, 2018 13th IEEE International Conference on. IEEE, pp. 263–270.

Lample, G., Zeghidour, N., Usunier, N., Bordes, A., Denoyer, L., et al., 2017. Fader networks: Manipulating images by sliding attributes. In: *Advances in Neural Information Processing Systems*. pp. 5969–5978.

Langner, O., Dotsch, R., Bijlstra, G., Wigboldus, D.H., Hawk, S.T., Van Knippenberg, A., 2010. Presentation and validation of the radboud faces database. *Cognition Emotion* 24 (8), 1377–1388.

Le, V., Brandt, J., Lin, Z., Bourdev, L., Huang, T.S., 2012. Interactive facial feature localization. In: *European Conference on Computer Vision*. Springer, pp. 679–692.

Ledig, C., Theis, L., Huszar, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., et al., 2017. Photo-realistic single image super-resolution using a generative adversarial network. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4681–4690.

Li, M., Zuo, W., Zhang, D., 2016. Deep Identity-aware Transfer of Facial Attributes, arXiv preprint arXiv:1610.05586.

Liu, M., Shan, S., Wang, R., Chen, X., 2014. Learning expressionlets on spatio-temporal manifold for dynamic facial expression recognition. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 1749–1756.

Lu, Y., Tai, Y.-W., Tang, C.-K., 2018. Attribute-guided face generation using conditional cylegan. In: *European Conference on Computer Vision*. Springer, pp. 293–308.

Odena, A., Olah, C., Shlens, J., 2017. Conditional image synthesis with auxiliary classifier gans. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR.org, pp. 2642–2651.

Parkhi, O.M., Vedaldi, A., Zisserman, A., et al., 2015. Deep face recognition. In: *BMVC*. p. 6.

Perarnau, G., van de Weijer, J., Raducanu, B., Álvarez, J.M., 2016. Invertible conditional gans for image editing, arXiv preprint arXiv:1611.06355.

Qiao, F., Yao, N., Jiao, Z., Li, Z., Chen, H., Wang, H., 2018. Geometry-Contrastive Generative Adversarial Network for Facial Expression Synthesis, arXiv preprint arXiv:1802.01822.

Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 234–241.

Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.-C., 2018. MobileNetV2: Inverted Residuals and Linear Bottlenecks, arXiv preprint arXiv:1801.04381.

Shen, W., Liu, R., 2017. Learning residual images for face attribute manipulation. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 1225–1233.

Shi, W., Caballero, J., Huszar, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z., 2016. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1874–1883.

Shrivastava, A., Pfister, T., Tuzel, O., Susskind, J., Wang, W., Webb, R., 2017. Learning from simulated and unsupervised images through adversarial training. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Vol. 3. p. 6.

Taigman, Y., Polyak, A., Wolf, L., 2016. Unsupervised cross-domain image generation, arXiv preprint arXiv:1611.02200.

Tariq, U., Yang, J., Huang, T.S., 2013. Maximum margin gmm learning for facial expression recognition. In: *Automatic Face and Gesture Recognition (FG)*, 2013 10th IEEE International Conference and Workshops on. IEEE, pp. 1–6.

- Vlasic, D., Brand, M., Pfister, H., Popović, J., 2005. Face transfer with multilinear models. *ACM Trans. Graphics (TOG)* 24 (3), 426–433.
- Yin, L., Wei, X., Sun, Y., Wang, J., Rosato, M.J., 2006. A 3D facial expression database for facial behavior research. In: *Automatic Face and Gesture Recognition, 2006. FGR 2006. 7th International Conference on*. IEEE, pp. 211–216.
- Zhang, F., Zhang, T., Mao, Q., Xu, C., 2018. Joint pose and expression modeling for facial expression recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3359–3368.
- Zhao, G., Huang, X., Taini, M., Li, S.Z., Pietikäinen, M., 2011. Facial expression recognition from near-infrared videos. *Image Vis. Comput.* 29 (9), 607–619.
- Zhu, X., Liu, Y., Li, J., Wan, T., Qin, Z., 2018. Emotion classification with data augmentation using generative adversarial networks. In: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, pp. 349–360.
- Zhu, J.-Y., Park, T., Isola, P., Efros, A.A., 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2223–2232.