# EXPRESSION CONDITIONAL GAN FOR FACIAL EXPRESSION-TO-EXPRESSION TRANSLATION

*Hao Tang[1], Wei Wang[2], Songsong Wu[3], Xinya Chen[4], Dan Xu[5], Nicu Sebe[1], Yan Yan[6]*

[1]University of Trento  [2]EPFL  [3]Nanjing University of Posts and Telecommunications
[4]Huazhong University of Science and Technology  [5]University of Oxford  [6]Texas State University

## ABSTRACT

In this paper, we focus on the facial expression translation task and propose a novel Expression Conditional GAN (EC-GAN) which can learn the mapping from one image domain to another one based on an additional expression attribute. The proposed ECGAN is a generic framework and is applicable to different expression generation tasks where specific facial expression can be easily controlled by the conditional attribute label. Besides, we introduce a novel face mask loss to reduce the influence of background changing. Moreover, we propose an entire framework for facial expression generation and recognition in the wild, which consists of two modules, i.e., generation and recognition. Finally, we evaluate our framework on several public face datasets in which the subjects have different races, illumination, occlusion, pose, color, content and background conditions. Even though these datasets are very diverse, both the qualitative and quantitative results demonstrate that our approach is able to generate facial expressions accurately and robustly.

***Index Terms***— Generative Adversarial Networks (GANs), Image-to-Image Translation, Facial Expression

## 1. INTRODUCTION

Recently, Generative Adversarial Networks (GANs) have shown to capture complex image data with numerous applications in computer vision and image processing. For example, Pix2pix [1] can translate an image from one domain to another one in a supervised way, i.e., the training image pairs are required. However, obtaining paired training data can be difficult and expensive in some cases as indicated in [2]. To tackle this limitation, Zhu et al. propose CycleGAN [2], in which the model can learn the mapping function from one domain to another one with unpaired training data. Similar ideas have been proposed in [3, 4, 5, 6]. Despite these efforts, facial expression translation remains a challenging task due to the fact that the expression changes are non-linear [7, 8].

To overcome the aforementioned challenging, we propose a novel Expression Conditional GAN (ECGAN) for facial expression translation based on CycleGAN [2]. ECGAN can generate faces with different emotions which are conditioned on the input expression attribute vector. Our work is inspired by IcGAN [9] which factorizes an input image into a latent representation and conditional information using the trained encoders. By changing the conditional information, the generator network combines the same latent representation and the changed conditional information to generate an image that satisfies the changed encoded constraints.

In this paper, we present another strategy in which the conditional attribute vector is concatenated with the image representation in the convolutional layers, as shown in Fig. 1. The conditional attribute is represented by a vector, which is used to distinguish each attribute from the others. In the attribute vector, only the element which corresponds to the label is set to 1 while the rest of them are set to 0. Then the vector is concatenated with the image embedding vector at the bottleneck which is a fully connected layer of generator $G_{X \rightarrow Y}$. In the generator $G_{Y \rightarrow X}$, we change the expression vector by swapping the corresponding two expressions. The conditional label can be used to guide the transformation from one expression to another one. For instance, as shown in Fig. 1, the anger label corresponds to an angry face with an open mouth. A correspondence between the anger label and the open mouth is built. During training time, GAN can learn this correspondence automatically. Thus, ECGAN can reshape a face (e.g., mouth) by adding the conditional vector.

Moreover, we introduce a novel face mask loss to reduce the influence of background changing similar to [10]. We also present a complete framework for facial expression translation and recognition in the wild. Our framework comprises of two modules, i.e., translation and recognition. ECGAN allows us how to map a face with neutral expression to the faces with other expressions (e.g., anger, disgust), and vice versa. We can explicitly control the expression of a face image via the conditional expression vector, which can be potentially useful in several applications, such as data augmentation and facial expression profiling. Then, we rely on a face recognition model to evaluate the generated images of ECGAN. Overall, our contribution is three-fold: (1) We propose EC-GAN, which allows us to generate and modify real images of faces conditioned on arbitrary facial expressions. (2) We propose a novel face mask loss for alleviating the influence of background changing. (3) We propose a new VGG score to evaluate the generated images by GAN models.
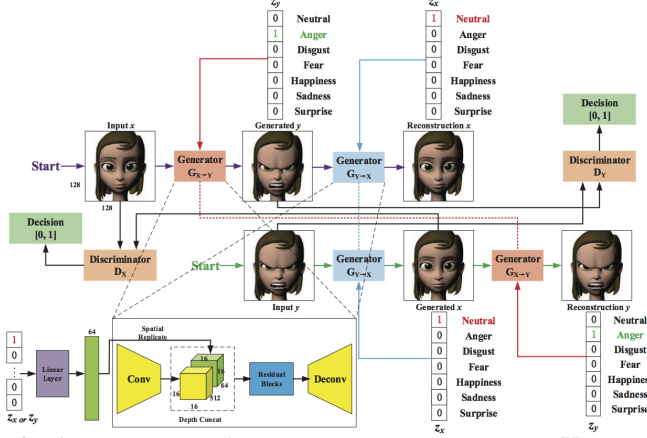
**Fig. 1**: Framework of the proposed ECGAN. Image $X$ can be converted to a new modified expression face image $Y$ guided by the facial expression attribute vector $z$. The expression attribute vector $z$ is concatenated with the image representation in the convolution layers.

## 2. RELATED WORK

**Generative Adversarial Networks (GANs)** [11] have achieved impressive performance on image generation tasks [12, 1, 2, 7, 13, 14]. Moreover, conditional GANs [15] are proposed to generate meaningful images that meet a certain requirement, where the conditioned label is employed to guide the image generation process. The conditioned labels can be discrete class labels [9], text descriptions [16, 17], object keypoints [18], human skeleton [14], semantic maps [19, 20] or reference images [1, 2]. The conditional models with images have tackled a lot of problems, e.g., image editing [9], text-to-image translation [16, 21], image-to-image translation [1] and video-to-video translation [22, 23].

**Image-to-Image Translation** learns a mapping function between different image domains using CNNs. Pix2pix [1] employs a conditional GAN to learn a mapping function from input to output images in a supervised way. Wang et al. [13] further propose Pix2pixHD model, which can turn semantic label maps into photo-realistic images or synthesizing portraits from face label maps. Similar ideas have also been applied to many other tasks, such as [24, 14, 25, 20]. However, these methods need to use the paired input-output data for training, which is not feasible for some applications. To overcome this limitation, Zhu et al. [2] propose CycleGAN, which learns the mappings between two different image domains with the unpaired data. Moreover, many other GAN variants are proposed to tackle the unpaired image-to-image translation task, such as [3, 4, 26, 27, 28, 5, 6, 29, 30, 31, 8, 32, 7].

**Face Editing.** Face analysis has a wide range of applications, such as face completion [33], hair modeling [34], aging [35], image-to-sketch translation [36, 37]. For example, Taigman et al. [4] propose Domain Transfer Network (DTN) for face-to-emoji translation task. Several other works [9, 31, 38] focus on human face attributes (e.g., bald, bangs, black hair, blond hair, eyeglasses, heavy makeup, male, mustache, pale skin)

translation. For instance, Larsen et al. [38] use a combination of Variational Autoencoder (VAE) and GAN to generate face samples with visual attribute vectors added to their latent representations. Shu et al. [39] present an end-to-end GAN that infers a face-specific disentangled representation of intrinsic face properties, including shape (i.e., normals), albedo, lighting, and an alpha matte. In this work, we focus on the arbitrary facial expression translation task with unpaired training data.

## 3. FORMULATION

GANs [11] are composed of two competing modules, i.e., a generator $G_{X \to Y}$ and a discriminator $D_Y$ (Where $X$ and $Y$ denote two different domains), which are iteratively trained competing against with each other in the manner of two-player minimax. CycleGAN [2] includes two mappings $G_{X \to Y}: X \to Y$ and $G_{Y \to X}: Y \to X$, and two adversarial discriminators $D_X$ and $D_Y$. The generator $G_{X \to Y}$ maps $X$ from the source domain to the target domain $Y$ and tries to fool the discriminator $D_Y$, whilst the $D_Y$ focuses on improving itself in order to be able to tell whether a sample is a generated sample or a real data sample. The similar to the generator $G_{Y \to X}$ and the discriminator $D_X$. More formally, let $x_i \in X$ and $y_j \in Y$ (For simplicity, we usually omit the subscript $i$ and $j$.) denote the training images in source and target image domain, respectively. We intent to learn a mapping function between $X$ domain and $Y$ domains with training data $\{x_i\}_{i=1}^N$ and $\{y_j\}_{j=1}^M$.

### 3.1. Objective Function

Our ECGAN objective contain several losses, we will introduce each of them, respectively.

**Adversarial loss.** We apply a least square loss [40] to stabilize our model during training. The least square loss is more stable than the negative log likelihood objective and more faster than Wasserstein GAN (WGAN) [41] to converge:

$$
\begin{aligned}
\mathcal{L}_{lsgan}(G_{X \to Y}, D_Y, X, Y) &= \mathbb{E}_{y \sim p_{\text{data}}(y)}[(D_Y(y) - 1)^2] \\
&+ \mathbb{E}_{x \sim p_{\text{data}}(x), z \sim p_z(z)}[D_Y(G_{X \to Y}(x, z))^2]
\end{aligned}
\tag{1}
$$

where $G_{X \to Y}$ tries to generate images $G_{X \to Y}(x, z)$ that look similar to images from domain $Y$, while $D_Y$ aims to distinguish between translated samples $G_{X \to Y}(x, z)$ and real samples $y$. $G_{X \to Y}$ aims to minimize this objective against an adversary $D_Y$ that tries to maximize it. We have a similar loss for generator $G_{Y \to X}$ and discriminator $D_X$ as well.

**Cycle Consistency Loss.** Note that CycleGAN [2] is different from Pix2pix [1] in the way that the training data in CycleGAN is unpaired. CycleGAN introduces the cycle consistency loss to enforce forward-backward consistency. The cycle consistency loss can be regarded as "pseudo" pairs of training data even though we do not have the corresponding data in the target domain which corresponds to the input data from the source domain. To include facial expression conditional constraint $z$ as part of the input to the generator and
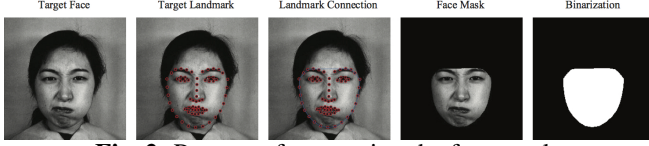
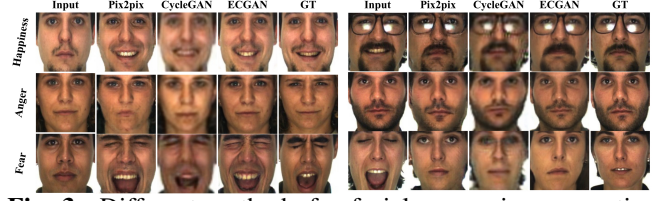**Fig. 2**: Process of computing the face mask.



**Fig. 3**: Different methods for facial expression generation (Left) and neutralization (Right). Form left to right: input, Pix2pix trained on paired data, CycleGAN, ECGAN and Ground Truth (GT). Note that images are cropped for visualization.

discriminator of CycleGAN, the loss of cycle consistency is reformulated as follows:

$$
\begin{aligned}
\mathcal{L}_{cyc}&(G_{X \to Y}, G_{Y \to X})\\
=&\mathbb{E}_{x \sim p_{\text{data}}(x), z \sim p_z(z)}[\|G_{Y \to X}(G_{X \to Y}(x, z)) - x\|_1]\\
+&\mathbb{E}_{y \sim p_{\text{data}}(y), z \sim p_z(z)}[\|G_{X \to Y}(G_{Y \to X}(y, z)) - y\|_1].
\end{aligned}
\tag{2}
$$

**Context Loss.** The pixel-wise MSE loss is used as a context loss in [42, 43]. However, since the pixel-wise MSE loss often lacks high-frequency content which results in perceptually unsatisfying solutions with overly smooth textures. Ledig et al. [44] introduces the VGG loss, which is closer to perceptual similarity. The formulation of the VGG loss as follows:

$$
\mathcal{L}_{content}^{VGG_{Y \to X}} = \frac{1}{W_{i,j}H_{i,j}} \sum_{w=1}^{W_{i,j}} \sum_{h=1}^{H_{i,j}} (\phi_{i,j}(X)_{w,h} - \phi_{i,j}G_{Y \to X}(y)_{w,h})^2,
\tag{3}
$$

where, $\phi_{i,j}$ indicate the feature map obtained by the $j$-th convolution before the $i$-th max-pooling layer within VGG net [45], $W_{i,j}$ and $H_{i,j}$ are the dimensions of the respective feature maps within the VGG network. Therefore, the final loss $\mathcal{L}_{content}^{VGG} = \mathcal{L}_{content}^{VGG_{Y \to X}} + \mathcal{L}_{content}^{VGG_{X \to Y}}$.

**Identity Preserving Loss.** To reinforce the identity of the face while converting, a face identity preserving loss [4] is adopted to preserve the identity.

$$
\mathcal{L}_{identity}(G_{X \to Y}, G_{Y \to X}) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\|G_{Y \to X}(x) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)}[\|G_{X \to Y}(y) - y\|_1]
\tag{4}
$$

In such way, generators will take into consideration the identity problem through the back-propagation of the identity loss.

**Face Mask Loss.** In order to eliminate the influence brought by background changes, we propose a novel loss that add a face mask $M$ to the $L_1$ loss such that the face is given larger weight than the background, as shown in Fig. 2. We apply OpenFace [46] to extract face landmark. The formulation of face mask is given as follows with $\odot$ as the pixel-wise multiplication:

$$
\mathcal{L}_{mask}^{Y \to X} = \|(G_{Y \to X}(G_{X \to Y}(x \odot M_x)) - x \odot M_x)\|_1,
\tag{5}
$$

face mask $M_x$ are set to 1 for foreground and 0 for background and applying a set of morphological operations such that it is able to approximately cover the whole face. Thus, $\mathcal{L}_{mask} = \mathcal{L}_{mask}^{Y \to X} + \mathcal{L}_{mask}^{X \to Y}$.

**Full Objective.** Consequently, the complete objective loss is:

$$
\begin{aligned}
\mathcal{L}(G_{X \to Y}, G_{Y \to X}, D_X, D_Y) =& \mathcal{L}_{cGAN} + \lambda_1 \mathcal{L}_{cyc}(G_{X \to Y}, G_{Y \to X}) +\\
& \lambda_2 \mathcal{L}_{content} + \lambda_3 \mathcal{L}_{identity} + \lambda_4 \mathcal{L}_{mask},
\end{aligned}
\tag{6}
$$

where $\lambda_1$, $\lambda_2$, $\lambda_3$ and $\lambda_4$ are parameters controlling the relative relation of objectives terms.

## 4. EXPERIMENTS

In this section, we first introduce the details of the employed datasets in our experiments, then we demonstrate the results and discussions of generation and recognition steps respectively.

**Datasets.** We employ several datasets to validate our model. These datasets contains faces with different races and they have different illumination, occlusion, pose conditions and backgrounds. See the supplementary materials for details.

**Setup.** We use the same training setups as CycleGAN [2]. Adam optimizer [47] with a batch size of 1 is used. The initial learning rate for Adam optimizer is 0.0002 and $\beta_1$ of Adam is 0.5. For fair comparisons, all models were trained for 200 epochs. Training and testing stages are conducted out on an Nvidia TITAN Xp GPU with 12GB memory.

**Competing Models.** We employ state-of-the-art image translation models, i.e., CycleGAN [2], Pix2pix [1] as our baselines. Note that Pix2pix [1] is trained on paired data. For a fair comparison, we implement both baselines using the same setups as our approach.

**Evaluation Metrics.** We provide both qualitative and quantitative results. Qualitatively, the images generated by different methods as shown in Fig. 3. Quantitatively, the expression recognition accuracy score is employed to evaluate whether the generated images wear the correct expressions. To this end, we propose a novel VGG Score which is similar to "FCN Score" in [1] and Inception Score [48] as the score of accuracy. The definition of the Inception Score is $exp(E_x[KL(p(y|x)\|p(y))])$, where $x$ is an image, $p(y|x)$ is the inferred class label probabilities given $x$ by the pre-trained Inception network and $p(y)$ is the marginal distribution over all images. The VGG score is defined as $E_x(p(y|x, z))$, where $z$ is the conditioning label. Overall, the differences between VGG score and Inception Score [48] are as follows. First, the VGG Score is calculated using a pre-trained VGG network, while the Inception Score is calculated using a pre-trained Inception network. Second, even though both the VGG Score and the Inception Score are defined to maximize inferred probabilities in order to guarantee that the generated images are meaningful, the difference is that the Inception Score includes an extra term to maximize the entropy of the marginal distributions to encourage the diversity of the gener-
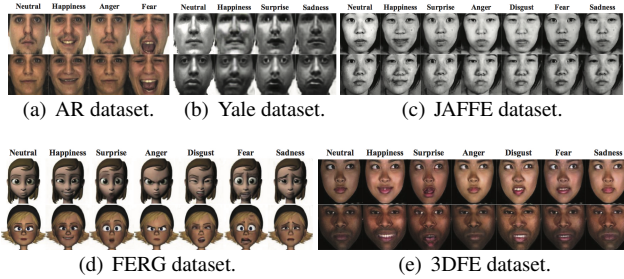
(a) AR dataset.    (b) Yale dataset.    (c) JAFFE dataset.

(d) FERG dataset.    (e) 3DFE dataset.

**Fig. 4**: Example results of ECGAN. The neutral expression is the input and the others expressions are the output. We can observe that even though the subjects in all datasets with significant differences, our method consistently generates high-quality images, which shows our method is very insensitive to changing skin color, posture, illumination or occlusion.

**Table 1**: AMT Score of different methods.

| Method | CycleGAN [2] | Pix2pix [1] | ECGAN (Ours) |
|---|---|---|---|
| AMT Score | 11.68 | **40.37** | 35.32 |

**Table 2**: VGG Score (%) of different methods.

| Method | Train Set | Test Set | VGG Score |
|---|---|---|---|
| baseline | original | original | 74.77 |
| CycleGAN [2] | +generated | original | 76.41 |
| CycleGAN [2] | original | generated | 77.78 |
| Pix2pix [1] | +generated | original | 82.63 |
| Pix2pix [1] | original | generated | **83.24** |
| ECGAN (Ours) | +generated | original | 78.13 |
| ECGAN (Ours) | original | generated | 80.32 |

ated images.

**Qualitative Evaluation.** Fig. 3 demonstrates the images generated by our method and the baselines. We can observe that the results generated by CycleGAN tend to be more blurry compared with Pix2pix and ECGAN. ECGAN adopts the proposed face mask loss to guide the generators to focus on the face regions. Though Pix2pix also generates images with competitive quality, the model can only be trained with paired data. In contrast, our ECGAN produces good quality images without the requirement of paired data. To exhaustively validate the superiority of our ECGAN, Fig. 4 provides more generation results. We can see that ECGAN generalizes well to the unseen data. We also observe that even though the subjects in all datasets have different races, poses, skin colors, illumination conditions and occlusions, our method consistently generates high-quality images. This demonstrates that our method is very robust.

**Quantitative Evaluation.** We follow [2] to conduct the "real vs fake" perceptual studies on Amazon Mechanical Turk (AMT) to assess the realism of the generated images. Results are shown in Table 1. We can see that the proposed method is significantly better than CycleGAN, but a little worse than Pix2pix. Moreover, we employ the expression recognition accuracy to evaluate the correctness of the generated expressions. The intuition is that if the generated images are realistic, then (i) the classifiers trained on both the real images and the generated images will be able to boost the accuracy of
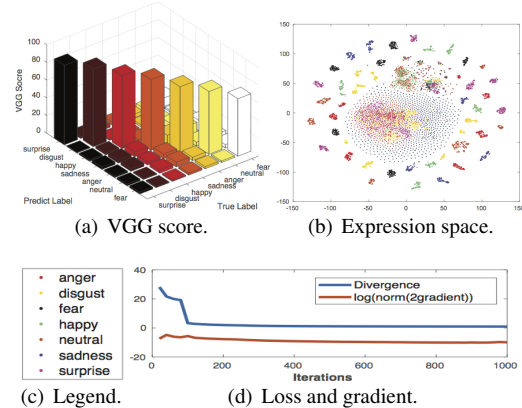


(a) VGG score.    (b) Expression space.

(c) Legend.    (d) Loss and gradient.

**Fig. 5**: (a) VGG score. (b-d) Feature space of the generated facial expressions. Each color represents a expression.

the real images (in this situation, the generated images work as augmented data.) (ii) the classifiers trained on real images will also be able to classify the synthesized image correctly. VGG [45] is adopted as deep feature extractor for our facial expression recognition task. Detailed recognition performance is reported in Table 2, when the model is trained only with the original training data, and tested on the testing data, the score is 74.77%. When the generated images are added to the training set as augmented data, the recognition score of the testing data is increased to 78.13%. Besides, to validate that our model can generate the correct expressions, we replaced the test set by the generated images and achieve 80.32% recognition rate (Fig. 5(a)), which demonstrates the effectiveness of our method since the generative images have a slight better performance than the ground truth images. Moreover, we also conduct the experiment of expression clustering to visulize the distribution of the generated images, t-SNE [49] is adopted to visualize the 4,096-D deep feature on a two dimensional space. Fig. 5(b-d) illustrates the deep feature space of the generated images as well as the evolving of the loss and gradient in the training stage. Note that the generated images with the same expressions are classified into the same clusters according to their representations in the deep feature space, which reveals that our ECGAN method can generate images with correct expressions.

## 5. CONCLUSION

We propose Expression Conditional GAN (ECGAN) for the facial expression generation task. The main technical contribution is the proposed Conditional CycleGAN which utilizes the expression label to guide the facial expression generation process. In ECGAN, the adversarial loss is modified to include a conditional expression feature vectors as parts of the inputs to the generator and discriminator networks. The expression attribute vector is utilized to represent the expression label. Experimental results demonstrate that our method not only presents compelling results but also achieves competitive results on facial expression recognition task.

Cisco, Inc for this research.

# 6. REFERENCES

[1] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017.

[2] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017.

[3] Ming-Yu Liu, Thomas Breuel, and Jan Kautz, "Unsupervised image-to-image translation networks," in *NIPS*, 2017.

[4] Yaniv Taigman, Adam Polyak, and Lior Wolf, "Unsupervised cross-domain image generation," in *ICLR*, 2017.

[5] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *ICML*, 2017.

[6] Zili Yi, Hao Zhang, Ping Tan Gong, et al., "Dualgan: Unsupervised dual learning for image-to-image translation," in *ICCV*, 2017.

[7] Hao Tang, Dan Xu, Nicu Sebe, and Yan Yan, "Attention-guided generative adversarial networks for unsupervised image-to-image translation," in *IJCNN*, 2019.

[8] Albert Pumarola, Antonio Agudo, Aleix M Martinez, Alberto Sanfeliu, and Francesc Moreno-Noguer, "Ganimation: Anatomically-aware facial animation from a single image," in *ECCV*, 2018.

[9] Guim Perarnau, Joost van de Weijer, Bogdan Raducanu, and Jose M Álvarez, "Invertible conditional gans for image editing," in *NIPS Workshop*, 2016.

[10] Xiaodan Liang, Hao Zhang, Liang Lin, and Eric Xing, "Generative semantic manipulation with mask-contrasting gan," in *ECCV*, 2018.

[11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial nets," in *NIPS*, 2014.

[12] Andrew Brock, Jeff Donahue, and Karen Simonyan, "Large scale gan training for high fidelity natural image synthesis," in *ICLR*, 2019.

[13] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro, "High-resolution image synthesis and semantic manipulation with conditional gans," in *CVPR*, 2018.

[14] Hao Tang, Wei Wang, Dan Xu, Yan Yan, and Nicu Sebe, "Gesturegan for hand gesture-to-gesture translation in the wild," in *ACM MM*, 2018.

[15] Mehdi Mirza and Simon Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[16] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov, "Generating images from captions with attention," in *ICLR*, 2015.

[17] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee, "Generative adversarial text-to-image synthesis," in *ICML*, 2016.

[18] Scott E Reed, Zeynep Akata, Santosh Mohan, Samuel Tenka, Bernt Schiele, and Honglak Lee, "Learning what and where to draw," in *NIPS*, 2016.

[19] Sangwoo Mo, Minsu Cho, and Jinwoo Shin, "Instagan: Instance-aware image-to-image translation," in *ICLR*, 2019.

[20] Hao Tang, Dan Xu, Nicu Sebe, Yanzhi Wang, Jason J Corso, and Yan Yan, "Multi-channel attention selection gan with cascaded semantic guidance for cross-view image translation," in *CVPR*, 2019.

[21] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas, "Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks," in *ICCV*, 2017.

[22] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro, "Video-to-video synthesis," in *NeurIPS*, 2018.

[23] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh, "Recycle-gan: Unsupervised video retargeting," in *ECCV*, 2018.

[24] Patsorn Sangkloy, Jingwan Lu, Chen Fang, Fisher Yu, and James Hays, "Scribbler: Controlling deep image synthesis with sketch and color," in *CVPR*, 2017.

[25] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman, "Toward multimodal image-to-image translation," in *NIPS*, 2017.

[26] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo, "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation," in *CVPR*, 2018.

[27] Sagie Benaim and Lior Wolf, "One-sided unsupervised domain mapping," in *NIPS*, 2017.

[28] Hao Tang, Dan Xu, Wei Wang, Yan Yan, and Nicu Sebe, "Dual generator generative adversarial networks for multi-domain image-to-image translation," in *ACCV*, 2018.

[29] Aaron Gokaslan, Vivek Ramanujan, Daniel Ritchie, Kwang In Kim, and James Tompkin, "Improving shape deformation in unsupervised image-to-image translation," in *ECCV*, 2018.

[30] Youssef A Mejjati, Christian Richardt, James Tompkin, Darren Cosker, and Kwang In Kim, "Unsupervised attention-guided image to image translation," in *NeurIPS*, 2018.

[31] Yongyi Lu, Yu-Wing Tai, and Chi-Keung Tang, "Conditional cyclegan for attribute guided face image generation," in *ECCV*, 2018.

[32] Huiwen Chang, Jingwan Lu, Fisher Yu, and Adam Finkelstein, "Paired-cyclegan: Asymmetric style transfer for applying and removing makeup," in *CVPR*, 2018.

[33] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang, "Generative face completion," in *CVPR*, 2017.

[34] Menglei Chai, Linjie Luo, Kalyan Sunkavalli, Nathan Carr, Sunil Hadap, and Kun Zhou, "High-quality hair modeling from a single portrait photo," *ACM TOG*, vol. 34, no. 6, pp. 204, 2015.

[35] Grigory Antipov, Moez Baccouche, and Jean-Luc Dugelay, "Face aging with conditional generative adversarial networks," in *ICIP*, 2017.

[36] Wengling Chen and James Hays, "Sketchygan: Towards diverse and realistic sketch to image synthesis," in *CVPR*, 2018.

[37] Hao Tang, Xinya Chen, Wei Wang, Dan Xu, Jason J Corso, Nicu Sebe, and Yan Yan, "Attribute-guided sketch generation," in *FG*, 2019.

[38] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther, "Autoencoding beyond pixels using a learned similarity metric," in *ICML*, 2015.

[39] Zhixin Shu, Ersin Yumer, Sunil Hadap, Kalyan Sunkavalli, Eli Shechtman, and Dimitris Samaras, "Neural face editing with intrinsic image disentangling," in *CVPR*, 2017.

[40] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley, "Least squares generative adversarial networks," in *ICCV*, 2017.

[41] Martin Arjovsky, Soumith Chintala, and Léon Bottou, "Wasserstein generative adversarial networks," in *ICML*, 2017.

[42] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang, "Image super-resolution using deep convolutional networks," *IEEE TPAMI*, vol. 38, no. 2, pp. 295–307, 2016.

[43] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *CVPR*, 2016.

[44] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al., "Photo-realistic single image super-resolution using a generative adversarial network," in *CVPR*, 2017.

[45] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.

[46] Brandon Amos, Bartosz Ludwiczuk, and Mahadev Satyanarayanan, "Openface: A general-purpose face recognition library with mobile applications," Tech. Rep., CMU-CS-16-118, CMU School of Computer Science, 2016.

[47] Diederik Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.

[48] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen, "Improved techniques for training gans," in *NIPS*, 2016.

[49] Laurens van der Maaten and Geoffrey Hinton, "Visualizing data using t-sne," *JMLR*, vol. 9, no. Nov, pp. 2579–2605, 2008.