

LI-NET: LARGE-POSE IDENTITY-PRESERVING FACE REENACTMENT NETWORK

Jin Liu^{*,†}, Peng Chen^{*,†}, Tao Liang^{*,†}, Zhaoxing Li^{*}, Cai Yu^{*,†}, Shuqiao Zou^{*,†}, Jiao Dai^{*}, Jizhong Han^{*}

^{*}Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

[†]School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

{liujin, chenpeng, liangtao0305, lizhaoxing, caiyu, zoushuqiao, daijiao, hanjizhong}@iie.ac.cn

ABSTRACT

Face reenactment is a challenging task, as it is difficult to maintain accurate expression, pose and identity simultaneously. Most existing methods directly apply driving facial landmarks to reenact source faces and ignore the intrinsic gap between two identities, resulting in the identity mismatch issue. Besides, they neglect the entanglement of expression and pose features when encoding driving faces, leading to inaccurate expressions and visual artifacts on large-pose reenacted faces. To address these problems, we propose a Large-pose Identity-preserving face reenactment network, LI-Net. Specifically, the Landmark Transformer is adopted to adjust driving landmark images, which aims to narrow the identity gap between driving and source landmark images. Then the Face Rotation Module and the Expression Enhancing Generator decouple the transformed landmark image into pose and expression features, and reenact those attributes separately to generate identity-preserving faces with accurate expressions and poses. Both qualitative and quantitative experimental results demonstrate the superiority of our method.

Index Terms— Face Reenactment, Image Synthesis, Generative Adversarial Network

1. INTRODUCTION

Given a source face and driving face, *face reenactment* aims to generate a reenacted face which is animated by the expression and pose of the driving face while preserving the identity of the source face. Various applications benefit from face reenactment, including telepresence, gaming and filmmaking. Recently, diverse face reenactment methods have emerged. However, existing methods suffer from identity mismatch and inaccurate expressions in large pose. In this work, we propose a Large-pose Identity-preserving face reenactment framework LI-Net to address the above problems.

Current face reenactment methods mainly rely on 3D models and GAN models. 3D-based methods reconstruct source faces using pre-defined 3DMM [1] and render reenacted faces on the image plane, which suffers from tedious

This research is supported in part by the National Key Research and Development Program of China under Grant No.2020AAA0140000.

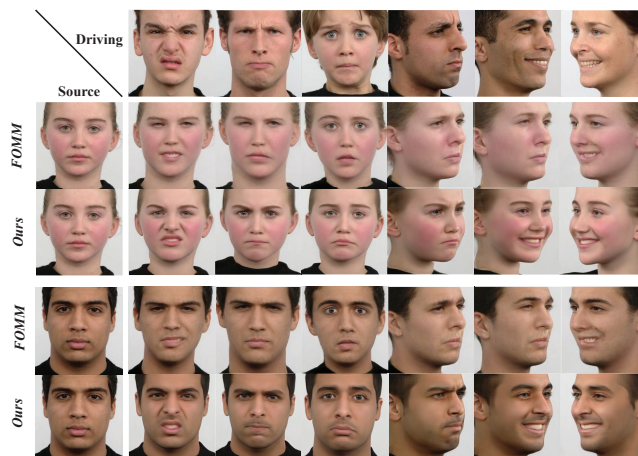


Fig. 1: The results of LI-Net and current SOTA method FOMM [15]. The first column and row are source and driving images. The other four rows are reenacted faces using methods notated on the left. Note we do not have identity mismatch or inaccurate expressions in large pose.

big-budget model construction procedure and high computation cost. For GAN-based methods, various approaches [21, 23] apply simple encoder-decoder network to reenact faces, restricted by predefined identities. Later, diverse many-to-many methods arise, which reenact any *unseen* identity using motion information from *unseen* driving faces. They either utilize large-scale data to improve generalization [15] or adopt image-warping to make full use of facial details in source faces [20]. Nonetheless, the above methods ignore the intrinsic gap between two identities and directly apply driving faces to reenact source faces, leading to identity mismatch problem. Besides, they reenact new expressions and poses simultaneously with the entanglement of these attributes, leading to inaccurate expressions and visual artifacts in large pose.

In detail, *identity mismatch* means the inability of face reenactment model to preserve the identity of the source face. We take the SOTA method FOMM [15] as an example. As shown in the first three columns in Fig. 1, we cannot tell the source identity from reenacted faces of FOMM whose facial contours are exactly the same as driving faces instead of source ones. Furthermore, when it comes to large-pose reen-

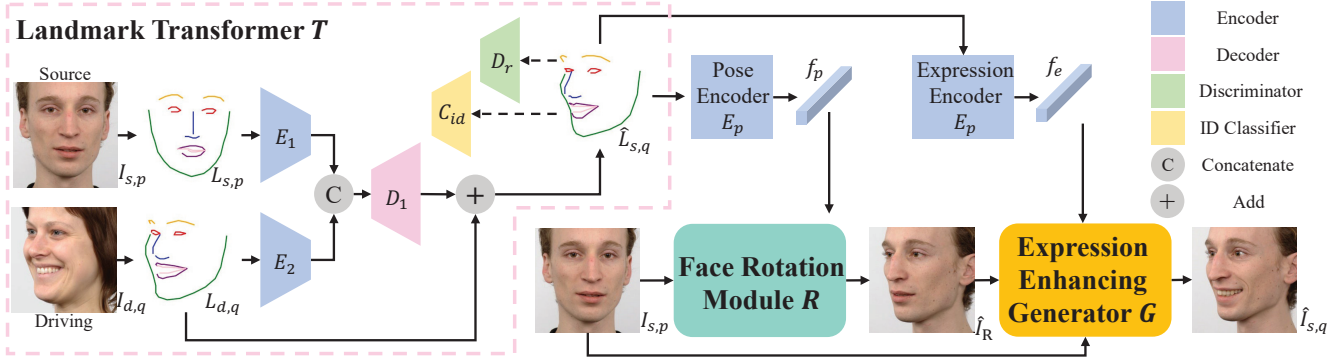


Fig. 2: Overview of the proposed LI-Net. First, the landmark transformer T (Section 3.1) takes $L_{s,p}$ and $L_{d,q}$ as inputs to generate transformed landmarks $\hat{L}_{s,q}$. The discriminator D_r and ID classifier C_{id} constrain the realism and identity of $\hat{L}_{s,q}$. Second, the face rotation module R (Section 3.2) will rotate source images according to pose information from $\hat{L}_{s,q}$. Finally, the expression enhancing generator G (Section 3.3) will enhance \hat{I}_R and add expressions to get reenacted faces $\hat{I}_{s,q}$.

actment, *inaccurate expressions* and *visual artifacts* appear. Fuzzy mouth shape and weird facial contour can be found in FOMM results, as shown in the last three columns of Fig. 1.

To address the above problems, we propose a Large-pose Identity-preserving face reenactment network called LI-Net. Specifically, given source and driving landmark images, the *Landmark Transformer* is utilized to transform driving landmarks to source identities. The *Face Rotation Module* then takes source face and transformed landmarks as input to generate rotated source faces with driving poses. Finally, source faces, rotated faces and transformed landmarks will be sent into *Expression Enhancing Generator* to get reenacted faces.

To the best of our knowledge, the proposed LI-Net is the first to perform many-to-many identity-preserving face reenactment while manages to maintain accurate identity, pose and expression simultaneously. Our contributions are as follows: (a) To solve the identity mismatch problem, the landmark transformer with explicit identity restrictions is proposed to transform driving landmarks to source identities. (b) We adopt the face rotation module with a novel data augmentation method and apply the expression enhancing generator to impose fine-grained pose and expression control over reenacted faces. (c) Both qualitative and quantitative experimental results indicate that the proposed method manages to reenact large-pose identity-preserving high-quality faces with accurate expressions and photo-realistic facial details.

2. RELATED WORK

Conditional Generative Adversarial Network. Diverse methods related to GAN [4] are proposed to achieve facial attribute editing or whole face image generation. cGAN [14] controls the attribute of faces using condition information vectors. Pix2Pix [8] takes face sketch photo as conditions to generate faces. MaskGAN [12] edits faces by mapping the segmentation mask to the target image. StyleGAN [9] applies transformation on latent codes to control face image styles using adaptive instance normalization (AdaIN) [7]. Unlike

these methods, LI-Net maps input driving faces to landmark images, which owns more scalability than vectors and is easy to obtain compared with segmentation masks.

Face Reenactment. Face reenactment mainly falls in two categories, 3D-based and GAN-based methods. 3D-based methods [18, 17] utilizes pre-defined 3DMM [1] to reconstruct shape, expression, texture and illumination of input faces and reenact faces by changing the corresponding parameters, which always reflects specific human races collected to construct 3DMM. For GAN-based methods, ReenactGAN [21] maps face pose information to a boundary latent space and reenact faces by a person-specific generator, restricted by the specific person identity. FReeNet [22] preprocesses landmarks and uses generator to apply many-to-many face reenactment, which can not transfer poses of driving faces, only expressions. X2Face [20] produces embedded faces from source faces and utilizes an encoder-decoder architecture. FOMM [15] predicts optical flow and occlusion masks to reenact faces using image warping. These methods suffer from identity mismatch and inaccurate expressions in large-pose faces. However, to solve the above problems, LI-Net uses a landmark transformer with explicit identity constraints and applies separate modules to impose fine-grained control over reenacted faces, regarding poses and expressions.

3. METHODOLOGY

The procedure of LI-Net is shown in Fig. 2. The landmark detector [2] first obtains landmark images $L_{s,p}$ and $L_{d,q}$ from faces $I_{s,p}$ and $I_{d,q}$. I and L denote faces and landmark images, respectively. The first subscript denotes identity while the second means motion information (expression and pose). Then the landmark transformer T transforms driving landmarks to source identities. Subsequently, based on a novel data augmentation method, the face rotation module R generates \hat{I}_R , containing pose information from $\hat{L}_{s,q}$. Finally, the expression enhancing generator G generates the final reenacted face $\hat{I}_{s,q}$, sharing the same pose and expression as $I_{d,q}$.

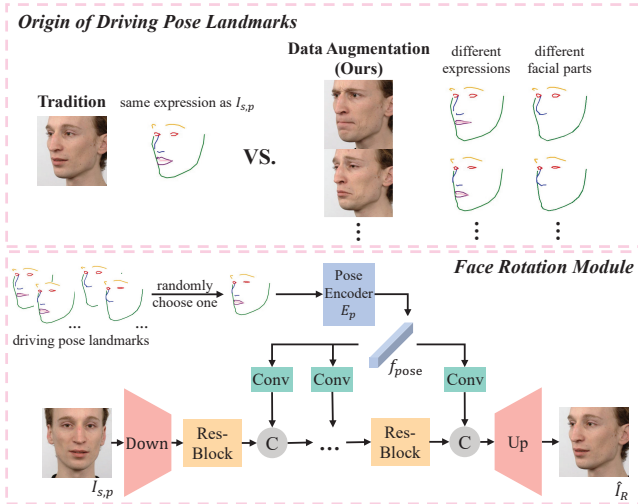


Fig. 3: Face Rotation module R and driving pose landmark images used in R . For a source image $I_{s,p}$ and its corresponding rotated image \hat{I}_R with driving pose, we extract pose features from face landmarks of same pose but with different expressions and facial parts. Each feature extracted will be sent into R to guide the same rotation process from $I_{s,p}$ to \hat{I}_R .

3.1. Landmark Transformer

To solve the identity mismatch problem, we propose a landmark transformer to transform driving landmark images to source face identities. As shown in Fig. 2, source face landmark $L_{s,p}$ and driving face landmark $L_{d,q}$ are sent to corresponding encoder E_1 and E_2 . Then the decoder D_1 transforms encoded features to predicted landmark shift, which will be added to $L_{d,q}$ and finally gets transformed reenacted landmarks $\hat{L}_{s,q}$. The landmark shift acts as fine adjustments on driving landmarks to give source identity information. Also, the identity classifier C_{id} and discriminator D_r guarantee the identity consistency and the realness of transformed landmarks. In general, we get the landmark transformer as $\hat{L}_{s,q} = T(L_{s,p}, L_{d,q})$, where the first and second input of T are source and driving landmarks, respectively. The overall loss function \mathcal{L}_T is defined as:

$$\mathcal{L}_T = \lambda_{L1} \mathcal{L}_{L1} + \lambda_{rec} \mathcal{L}_{rec} + \lambda_{cycle} \mathcal{L}_{cycle} + \lambda_{id} \mathcal{L}_{id} + \lambda_D \mathcal{L}_D, \quad (1)$$

where all λ are weights of various loss terms. We define \mathcal{L}_{L1} as pixel-wise L1 loss between transformed landmark $\hat{L}_{s,q}$ and the corresponding ground truth landmark $L_{s,q}$.

Reconstruction Loss. The landmark transformer should not only transform landmarks between two different identities, but also reconstruct the driving landmark when receiving different landmarks with same identity. which is constrained by:

$$\mathcal{L}_{rec} = \|T(L_{s,p}, L_{s,q}) - L_{s,q}\|_1. \quad (2)$$

Cycle Loss. Inspired by [23], \mathcal{L}_{cycle} are added in order to guarantee that the transformed landmark $\hat{L}_{s,q}$ can be trans-

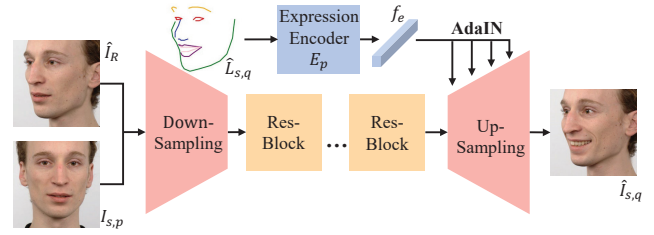


Fig. 4: Details of expression enhancing generator G . Given the rotated face \hat{I}_R , the source face $\hat{I}_{s,p}$ and the expression feature extracted from $\hat{L}_{s,q}$, G enhances \hat{I}_R and reenacts new expressions, generating $\hat{I}_{s,q}$.

formed back again to improve the fidelity, as shown below:

$$\mathcal{L}_{cycle} = \|T(L_{d,p}, T(L_{s,p}, L_{d,q})) - L_{d,q}\|_1. \quad (3)$$

Identity Loss. We introduce an identity classifier C_{id} to extract id features and predict id labels. During training, we use the cross entropy loss between predicted and ground truth identity label to constrain the classifier. The identity loss term \mathcal{L}_{id} in Equation 1 is defined as the mean absolute error between identity features of $L_{s,p}$ and $\hat{L}_{s,q}$.

Adversarial Loss. Regarding the landmark transformer T as a generator, we introduce the discriminator D_r to judge the realness of transformed landmarks, which is constrained by:

$$\mathcal{L}_D = \log D_r(L_{s,q}) + \log(1 - D_r(\hat{L}_{s,q})), \quad (4)$$

where $L_{s,q}$ denotes the ground truth corresponding to $\hat{L}_{s,q}$.

3.2. Face Rotation Module

The face rotation module R is trained to rotate the source faces to the pose shown in $\hat{L}_{s,q}$. It is ill-posed and difficult to directly infer a large part of the face from source face, *e.g.* the left face of source person in Fig. 2. Therefore, the ultimate goal of R is not to generate photo-realistic rotated faces, but to produce a contour-consistent face image with driving pose to focus on new pose generation and reduce the burden of the entire reenacting process, *i.e.* from $I_{s,p}$ to $\hat{I}_{s,q}$.

Traditionally, when rotating the source face to a driving pose, we choose driving landmark images with driving pose that shares the same identity and expression with source face, as shown in the upper left of Fig. 3. However, for the driving landmarks of face rotation task, we argue that only the face contour plays the most important role to decide the direction of faces while expressions and other facial parts can be relatively ignored. Hence, we propose a data augmentation method as shown in upper right of Fig. 3. We choose landmark images with driving pose and source identity but with different expressions and facial parts. When training, we randomly choose one from those and extract the pose feature, which will be sent into R along with source face $I_{s,p}$ to generate \hat{I}_R , as shown in the bottom half of Fig. 3. In general, the loss terms are as follows:

$$\mathcal{L}_R = \lambda_{diff} \mathcal{L}_{diff} + \lambda_{GAN} \mathcal{L}_{GAN} + \lambda_{pose} \mathcal{L}_{pose}, \quad (5)$$

where λ are weights of loss terms. \mathcal{L}_{diff} minimizes the difference between rotated image \hat{I}_R and its ground truth I_R .

GAN Loss. Taking the rotation module R as a generator, we apply the discriminator loss function based on LSGAN [13]:

$$\mathcal{L}_{GAN} = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}))^2] + \mathbb{E}_{\mathbf{z} \sim p_{\text{data}}(\mathbf{z})} [(D(R(\mathbf{z})) - 1)^2], \quad (6)$$

where x denotes real face image data space and z denotes the input data space of the face rotation module R .

Pose Prediction Loss. To improve the accuracy of face pose synthesis, we employ another discriminator D_p to predict the pose of a given landmark, which is constrained by:

$$\mathcal{L}_{D_p} = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D_p(\mathbf{x}) - p)^2], \quad (7)$$

$$\mathcal{L}_{pose} = \mathbb{E}_{\mathbf{z} \sim p_{\text{data}}(\mathbf{z})} [(D_p(R(\mathbf{z})) - p)^2],$$

where p denotes the driving pose vector. x and z follow the same meaning in GAN Loss. \mathcal{L}_{D_p} constrains D_p to predict the correct pose while \mathcal{L}_{pose} pushes module R to generate faces with driving pose.

3.3. Expression Enhancing Generator

As shown in Fig. 4, given the rotated source face \hat{I}_R and the transformed landmark $\hat{L}_{s,q}$, the expression enhancing generator G aims to enhance detailed information and reenact expression contained in $\hat{L}_{s,q}$ to finally get reenacted face image $\hat{I}_{s,q}$. An expression encoder E_e is utilized to encode expression information, which is later injected in residual blocks by AdaIN [7] layers. During training, the expression enhancing generator G simply changes the expression while preserving the identity and driving pose. The total loss \mathcal{L}_G is defined as:

$$\mathcal{L}_G = \lambda_{pix} \mathcal{L}_{pix} + \lambda_{per} \mathcal{L}_{per} + \lambda_{adv} \mathcal{L}_{adv}, \quad (8)$$

where λ are weights of various loss functions. We use \mathcal{L}_{pix} as pixel-wise L1 loss to minimize the pixel difference between the rotated image $\hat{I}_{s,q}$ and the ground truth image $I_{s,q}$. The term \mathcal{L}_{adv} is used to improve the realness of reenacted image $\hat{I}_{s,q}$ in an adversarial way using traditional GAN loss.

Perceptual Loss. \mathcal{L}_{per} is the perceptual loss for minimizing the semantic difference of images in feature level:

$$\mathcal{L}_{per} = \sum_{l \in \Phi} \left\| \phi_l(\hat{I}_{s,q}) - \phi_l(I_{s,q}) \right\| + \sum_{l \in \Psi} \left\| \psi_l(\hat{I}_{s,q}) - \psi_l(I_{s,q}) \right\|, \quad (9)$$

where Φ and Ψ are collections of convolution layers from the perceptual networks while Φ_l and Ψ_l are activation outputs from l -th layers. In order to constrain that reenacted images $\hat{I}_{s,q}$ share same semantic and identity information as $I_{s,q}$, two perceptual networks are utilized, which are VGG-19 [16] for image classification and VGGFace [3] for face verification.

4. EXPERIMENTS

4.1. Experimental Settings

Datasets. We use Radboud Faces Database (RaFD) [11] and Multi-PIE [5] Database for experiment. RaFD has 8040 im-



Fig. 5: Qualitative results on Multi-PIE dataset.

ages of 67 persons displaying 8 emotional expressions, each in five different angles. Faces are cropped to 256×256 and landmarks are detected by [2]. Faces of 55 persons are randomly selected as training set and the others test set, which means no identity overlap between two sets. Multi-PIE contains more than 750,000 images of 337 subjects, captured under 15 view points and 19 illumination conditions in four recording sessions. We use a subset of 75,000 images (2 illumination conditions) to reduce training time.

Implementation Details. Our method is implemented using PyTorch 1.5.1 on Ubuntu 18.04 with a Tesla V100 GPU. Three modules in LI-Net are trained separately using the Adam optimizer [10]. Landmark transformer T is trained for 2,000 epochs with batch size 128 and starting learning rate $1e^{-5}$. We train face rotation module R and expression enhancing generator G for 500 epochs with batch size 32. The initial learning rate is $2e^{-4}$ and decays by ten every 100 epochs. The λ unifies losses on the same order of magnitude. All three modules are trained separately and detailed training schemes are described in Section 3.

4.2. Experimental Results

Methods. We choose X2Face [20], Pix2Pix [8], FOMM [15] and FReeNet [22] as comparison methods. X2Face uses image warping to warp generated embedded face according to extracted driving expression vectors. For Pix2Pix, we follow typical image translation training scheme and concatenate source faces and driving landmark images as input. The SOTA method FOMM decouples appearance and motion information and predicts optical flow and occlusion masks to generate reenacted faces. FReeNet [22] pre-processes landmarks and directly produces reenacted faces. All methods are fine-tuned using RaFD based on their pre-trained models.

Quantitative Results. We compare results of face reenactment quantitatively on RaFD dataset using two metrics, structured similarity index (SSIM) [19] and Fréchet inception distance (FID) [6]. Higher SSIM and lower FID mean better reenactment results. During the experiment, for each identity in test set which serves as source faces, 400 driving landmarks from test set are randomly selected to generate reenacted faces

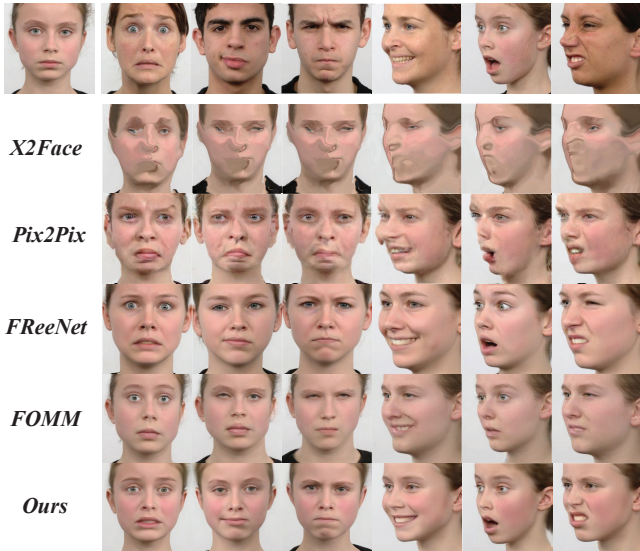


Fig. 6: Qualitative results between LI-Net and other methods on RaFD.

Table 1: Metric evaluation results between other methods and LI-Net on the RaFD dataset.

Model	SSIM \uparrow	FID \downarrow
X2Face [20]	0.61	128.43
Pix2Pix [8]	0.64	107.98
FOMM [15]	0.66	76.90
FReeNet [22]	<u>0.68</u>	83.31
Ours	0.73	<u>80.45</u>

(4,800 images totally). In this way, faces of various unseen identities, diverse expressions and poses are covered.

Table 1 shows the details of quantitative evaluation, demonstrating the effectiveness of our method. Note that for FReeNet, we do not copy the original paper evaluation results, because FReeNet only adopts partial RaFD and applies test identities to training, while we use full RaFD and different split method which means no identity overlap between training and test. LI-Net gets the best SSIM score and competitive FID score compared to FOMM. It is reasonable because FID judges both variety and reality of the image. FOMM may generate various faces given that diverse landmark images are taken as guidance for one identity, unlike LI-Net. Generated face details can be seen in Fig. 6.

Qualitative Results. Fig. 5 shows experimental results using Multi-PIE dataset. LI-Net can well transfer driving expressions to source persons while simultaneously maintain the correct pose and identity information. Fig. 6 shows reenacted faces between LI-Net and other methods. X2Face does not employ image realness discriminator and can not well grab face texture details simply using embedded faces and pose codes. Pix2Pix has large artifacts because the driving landmarks are used only once as the input, thus losing control to the reenacted faces when the net goes deeper. FReeNet and FOMM suffer from identity mismatch when reenacting un-

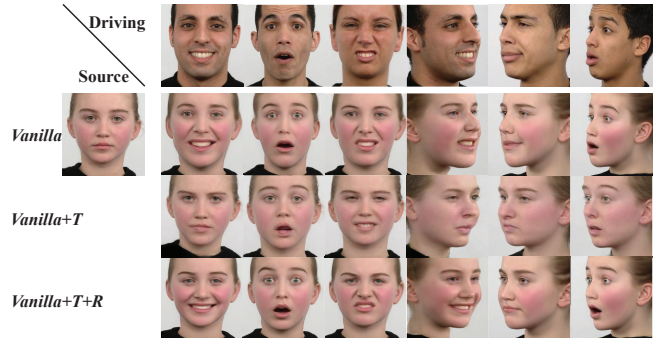


Fig. 7: Qualitative results of ablation studies on RaFD.

Table 2: Metric evaluation results of LI-Net with different components on the RaFD dataset.

Model	SSIM \uparrow	FID \downarrow
Vanilla	0.67	<u>83.30</u>
Vanilla +T	<u>0.68</u>	89.98
Vanilla +T+R	0.73	80.45

seen subjects, the inaccurate expression and artifacts in large pose, since they do not impose identity constraints to driving landmarks or explicitly deal with large pose situation to alleviate the inaccurate expressions or visual artifacts. Overall, LI-Net achieves the best performance in image quality while successfully preserving source identity and maintaining accurate expressions in large pose.

Ablation Studies. To demonstrate the effectiveness of each module detailedly, we conduct ablation studies on three different settings using RaFD dataset: (a) *vanilla*: we treat the entire process as image-to-image translation task and only utilize the expression enhancing generator G . (b) *vanilla+T*: based on (a), we add the landmark transformer T . (c) *vanilla+T+R*: our full proposed methods, *i.e.* LI-Net.

Table 2 displays the metric evaluation results of ablation studies which shows that each component improves the SSIM score. T successfully transforms the landmarks to the source identity. R separates rotation and expression enhancing, thus improving the accuracy of facial expressions. For FID score, T solves the identity mismatch problem but inevitably lacks the diversity in reenacted faces, which slightly gives a negative effect. When adding module R , photo-realistic accurate face images are generated with more facial details, greatly improving the FID. Qualitative results are more intuitive. As shown in Fig. 7, the face identity of vanilla results are heavily dominated by driving faces. The locations of eyes, nose, mouth between driving faces and reenacted faces are exactly the same, which is unrealistic. *Vanilla+T* generates accurate source identity and maintains source face contour. However, it causes inaccurate fuzzy expressions and visual artifacts in large pose when driven by the transformed landmark. When adding R , we get faces with less artifacts, more accurate expressions and face contours in large pose which can be found between *Vanilla+T* and *Vanilla+T+R* results. Overall, we see a significant contribution of each component.

5. CONCLUSION

We present a Large-pose Identity-preserving face reenactment network called LI-Net to solve the identity mismatch problem and to generate accurate expressions with no artifacts in large-pose reenacted faces. We propose the Landmark Transformer to transform driving landmark images and adopt separate modules to generate poses and expressions individually, leading to fine-grained control over the reenacted faces. Experimental results have demonstrated that our work achieves a good performance in large-pose identity-preserving face reenactment. In future work, we plan to extend our framework to handle faces in the wild with complex backgrounds and arbitrary expressions.

6. REFERENCES

- [1] V. Blanz and T. Vetter, "A morphable model for the synthesis of 3d faces," in *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999, pp. 187–194.
- [2] 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 *International Conference on Computer Vision*, 2017.
- [3] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vggface2: A dataset for recognising faces across pose and age," in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 67–74.
- [4] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems 27*, 2014, pp. 2672–2680.
- [5] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker, "Multi-pie," *Image and Vision Computing*, vol. 28, no. 5, pp. 807–813, 2010.
- [6] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," in *Advances in neural information processing systems*, 2017, pp. 6626–6637.
- [7] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1501–1510.
- [8] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *CVPR*, 2017.
- [9] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 4401–4410.
- [10] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [11] O. Langner, R. Dotsch, G. Bijlstra, D. H. Wigboldus, S. T. Hawk, and A. Van Knippenberg, "Presentation and validation of the radboud faces database," *Cognition and emotion*, vol. 24, no. 8, pp. 1377–1388, 2010.
- [12] C.-H. Lee, Z. Liu, L. Wu, and P. Luo, "Maskgan: Towards diverse and interactive facial image manipulation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [13] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2794–2802.
- [14] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [15] A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe, "First order motion model for image animation," in *Advances in Neural Information Processing Systems*, 2019, pp. 7137–7147.
- [16] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [17] S. Suwajanakorn, S. M. Seitz, and I. Kemelmacher-Shlizerman, "Synthesizing obama: learning lip sync from audio," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 1–13, 2017.
- [18] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Nießner, "Face2face: Real-time face capture and reenactment of rgb videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2387–2395.
- [19] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [20] O. Wiles, A. Sophia Koepke, and A. Zisserman, "X2face: A network for controlling face generation using images, audio, and pose codes," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 670–686.
- [21] W. Wu, Y. Zhang, C. Li, C. Qian, and C. Change Loy, "Reenactgan: Learning to reenact faces via boundary transfer," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 603–619.
- [22] J. Zhang, X. Zeng, M. Wang, Y. Pan, L. Liu, Y. Liu, Y. Ding, and C. Fan, "Freenet: Multi-identity face reenactment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5326–5335.
- [23] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.