# Make a Face: Towards Arbitrary High Fidelity Face Manipulation

Shengju Qian[1]* Kwan-Yee Lin[5] Wayne Wu[2,5] Yangxiaokang Liu[3] Quan Wang[5]
Fumin Shen[3] Chen Qian[5] Ran He[4]
[1]The Chinese University of Hong Kong [2]Tsinghua University
[3]University of Electronic Science and Technology of China [4]NLPR, CASIA
[5]SenseTime Research

Input                          Face Manipulation Results of Our Model



Figure 1: Face manipulation results on in-the-wild samples via transferring knowledge learned from the CelebA dataset. The first column shows input images and the remainders are images generated by AF-VAE with target expression/rotation boundary maps as the condition. Note that the model is fine-tuned with movie clip frames from YouTube of $256 \times 256$ resolution. All the generated poses are unseen before.

## Abstract

*Recent studies have shown remarkable success in face manipulation task with the advance of GANs and VAEs paradigms, but the outputs are sometimes limited to low-resolution and lack of diversity.*

*In this work, we propose Additive Focal Variational Auto-encoder (AF-VAE), a novel approach that can arbitrarily manipulate high-resolution face images using a simple yet effective model and only weak supervision of reconstruction and KL divergence losses. First, a novel additive Gaussian Mixture assumption is introduced with an unsupervised clustering mechanism in the structural latent space, which endows better disentanglement and boosts multi-modal representation with external memory. Second, to improve the perceptual quality of synthesized results, two simple strategies in architecture design are further tailored and discussed on the behavior of Human Visual System (HVS) for the first time, allowing for fine control over the model complexity and sample quality. Human opinion studies and new state-of-the-art Inception Score (IS) / Fréchet Inception Distance (FID) demonstrate the superiority of our approach over existing algorithms, advancing both the fidelity and extremity of face manipulation task.*

---

[1]Work done during an internship at SenseTime Research.

# 1. Introduction

Automatically manipulating facial expressions and head poses from a single image is a challenging open-end conditional generation task. Faithful photo-realistic face manipulation finds a wide range of applications in industry, such as film production, face analysis, and photography technologies. With the flourish of generative models, the state of this task has advanced dramatically in recent years at the forefront of efforts to generating diverse and photo-realistic results. Nevertheless, it is highly challenging for a generative model to learn a compact representation of intrinsic face properties to synthesize face images with *high-fidelity* or *large facial expression/poses*, due to the ill-posed nature of lacking paired training data.

Current state-of-the-art face manipulation approaches [32, 22, 9, 37] mainly benefit from the advancement of generative adversarial networks (GANs). To tackle the two bottlenecks mentioned above, vast algorithms focus on sophisticated modifications to loss term or generator architecture, with the injection of different facial attribute information [54, 37, 22, 48, 50, 34, 49, 27]. Other works focus on the design of task-specific training procedure [54, 16, 48, 55, 7, 2, 28, 29]. Nevertheless, successful generation of plausible samples on extreme face geomorphing and complex uncontrolled datasets remain elusive goals for these methods due to their unstable training procedures and environmental constraints. The current state-of-the-art [32] in RaFD [26] face expression synthesis achieves a Fréchet Inception Distance (FID) [15] of $34$, still leaving a large gap towards real data even in a controlled environment.

In this work, we set out to close the gap in fidelity and extremity between facial expressions/rotations generated by current state-of-the-arts and real-world face images with a simple yet effective framework. We explore the conditional variational auto-encoder (C-VAE) formalism [21, 41] for face manipulation task. It is intuitive to adopt C-VAE by taking advantage of its nice manifold representation and stable training mechanism. Nonetheless, tailoring the "vanilla" C-VAE to face manipulation task is *non-trivial*: (1) The diversity of synthesized outputs will be sacrificed since the latent distribution is commonly assumed to be a unit Gaussian. However, as visualized in Fig. 2 (b), a complex latent representation that needs to describe factors like ages, complexions, luminance, and poses would break the Gaussian assumption due to its insufficiency. (2) Facial expressions are more fine-grained than other sources like landscapes, digits, and animals. Thus the property and common architecture of VAE could not satisfy the requisites for maintaining facial details on high-resolution images.

As a solution to such problems, we propose a novel Additive focal variational auto-encoder (AF-VAE) framework. By applying a light-weight geometry-guidance to explicitly disentangle facial appearance and structure in latent space, we encourage the latent code to be separated into a pose-invariant appearance representation and a structure representation, thereby preserving the appearance and structure information under geomorphing. To tackle the issue of diversity, a novel additive memory module that bridges unsupervised clustering mechanism with Gaussian Mixture prior in the structured latent space is introduced to the framework, endowing the ability of multi-model facial expression/rotations generation.

To further improve the perceptual quality of synthesized results, we also discover two simple yet effective strategies in model design and characterize them empirically with the behavior of Human Visual System (HVS). Leveraging insights from this empirical analysis, we demonstrate that a simple injection of these strategies can easily improve the perceptual quality of synthesized results. Our model can steadily manipulate photo-realistic facial expression and face rotations at $256 \times 256$ resolution under uncontrolled settings. The proposed AM-VAE improves the state-of-the-art Fréchet Inception Distance (FID) and Inception Score (IS) on uncontrolled CelebA [26] from $71.3$ and $1.065$ to $36.82$ and $2.15$.

We conduct comparisons with current state-of-the-art face manipulation algorithms [9, 32, 37, 46] and show that our approach outperforms these methods regarding both quantitative and qualitative evaluation. Extensive self-evaluation experiments further demonstrate the effectiveness of proposed components.

# 2. Related Work

**Face Manipulation.** In the literature of face image manipulation, besides classic mass-spring models and 2D/3D morphing methods [42, 43, 44], recently, significant progress has been achieved by leveraging the power of Generative Adversarial Networks (GANs) [32, 22, 9, 37, 53, 51] for photo-realistic synthesis results. To improve the robustness and diversity of GANs, tweaks on various aspects are explored. For example, StarGAN [9] exploits cycle consistency to preserve key attributes between the source and target images. GANimation [32] takes a step further, utilizing denser AUs prior vector to increase the diversity by altering the magnitude of each AUs. Attention mechanism in generator architecture is also used to mask out irrelevant facial regions, enforcing the network to only synthesize in-region texture. FaceID-GAN [37] introduces a three-player adversarial scheme and utilizes 3DMM [3] to better preserve facial property. CAPG-GAN [16] trains a couple-agent discriminator to constrain both distributions of pose and facial structure.

Since GANs face challenges in brittle training procedure and sampling diversity, some works take advantages from variational auto-encoder (VAE) paradigm and its variants

with an exploration of disentangling intrinsic facial properties in latent space. For example, Neural Face Editing [40] and Deforming Autoencoders [39] utilize graphics rendering elements such as UV maps, albedo, and shading to decouple the latent representation. However, the main disadvantages of these VAE-based methods are the blurry synthesized results caused by injected element-wise divergence measurement and imperfect network architecture. CVAE-GAN [1] combines VAE and GAN into a framework with an asymmetric training loss and fine-grained category label. These methods lack an interface for users to manipulate facial expression arbitrarily since they edit the results through manifold traversal. Notably, the capacity constraint of VAE is also discussed in other tasks like image caption, in [45, 33], an additive Gaussian encoding space is proposed to provide a more diverse and accurate caption result. Motivated by those prior works, we exploit AF-VAE, which can compensate for the drawbacks of VAEs by providing a conditional geometry-related additive memory prior as well as two light-weight network design strategies to the framework.

**High Fidelity Image Synthesis.** Recently, generating high-resolution samples with fine details and realistic textures has become a trend in image synthesis task. For instance, pix2pixHD [47] introduces a coarse-to-fine generator, a multi-scale discriminator and a feature matching loss on the basis of pix2pix [20]. It could translate facial edge to photo-realistic face image under $1024 \times 1024$ resolution. While it requires paired training data and thereby cannot be generalized to arbitrary edges. [8] utilizes progressive GANs to generate $512 \times 512$ face images under the control of discrete one-hot attributes. IntroVAE [17] proposes to jointly train its inference and generator in an introspective way, achieving $1024 \times 1024$ resolution on reconstruction. BigGAN [4] modifies a regularization scheme and a sampling technique to class-conditional GANs to achieve 512 resolution on ImageNet. While all these frameworks are only capable of uncontrolled generation or discrete attribute editing, our framework focuses on challenging high fidelity face manipulation without leveraging paired training data.

## 3. Method

In this section, we explore methods for manipulating facial expression/rotation and scaling up modal training to reap the benefits of detailed architecture design. We first clarify our notations here. Given a face image $x$ from a dataset $X$, the goal of our task is to learn a mapping $\mathcal{G}$ to transfer $x$ to an output image $\tilde{x}$, conditioned on a target facial structural information $c$. The intrinsic challenge lying behind this task is to preserve facial *appearance* under the high-fidelity setting and dramatic *structural* changes. In this light, we decouple the mapping $\mathcal{G}$ into $\phi_{app}$ and $\mu_{str}$, which are expected to learn the pose-invariant appearance

representation $z = \phi_{app}(x, c)$ and rational structure representation $y = \mu_{str}(c)$ respectively.

In this way, the conditional variational auto-encoder (C-VAE) can be tailored as a baseline for its capability of disentangling $z$ from $y$ by maximizing the lower bound on the conditional data-log-likelihood $p(x|y)$, *i.e.*,

$$\log p(x|y) \geq \mathbb{E}_q[\log p(x|z, y)] - D_{KL}[q(z|x, y), p(z|y)],$$
(1)

where $q(z|x, y)$ is an approximate distribution of posterior $p(z|y)$. In particular, $q(z|x, y)$ and $p(x|z, y)$ are the *encoder* and *decoder* respectively. The model is typically trained with the following stochastic objective:

$$\mathcal{L}(x, \phi, \varphi) = -\frac{1}{N} \sum_{i=1}^{N} \log p_\varphi(x^i|z^i, y^i) +$$
(2)

$$D_{KL}[q_\phi(z|x, y), p(z|y)], s.t. \forall i \ z^i \sim q_\phi(z|x, y)$$

where $\phi$ and $\varphi$, the parameters for the encoder and decoder, are learned with reparameterization trick [21]. $q_\phi(z|x, y)$ is typically restricted to be a distribution over $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as in [21].

However, there are three problems to discuss under the formulation of problems: (1)How to accurately disentangle the latent space to guarantee that $z$ is complementary to the structure prior $y$, with nothing but appearance? (2) The KL divergence term facilitates the structure of the learnt latent space to get close to the prior $p(z|y)$ over $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Does this choice satisfy the encoder's modeling of the diversity of the face samples? (3) Does the perceptual quality of synthesized results look appealing to viewers with only the reconstruction and KL divergence losses?

We concern the optimality of these choices on face manipulation and explore better solutions in the following sections. An overview of our framework is illustrated in Fig. 2(a).

### 3.1. Geometry-Guided Disentanglement

In order to arbitrarily manipulate a face image, sufficient geometry information should be provided in $c$. Typical choices of manipulated condition include facial landmarks, 3DMM parameters and masks. On the basis of sparse landmarks information as used in Pix2pixHD [47], we use a off-line interpolation process in [52] to obtain boundary maps. This off-line process is formulated as $c = F_\omega(x)$. Then the latent structure representation $y$ could be obtained by encoding $c$ with an encoder $E_\mu$.

Without any semantic assumption, the framework can only make sure the latent code $z$ drawn from $x$ is invariant to structure $y$. To this end, we leverage the distilled geometry information $y$ to disentangle the latent space *explicitly*. Concretely, $y$ is concatenated with the inferred appearance representation $z$. Then, the concatenated representation is forwarded to decoder $D_\varphi$. The skip-connections between
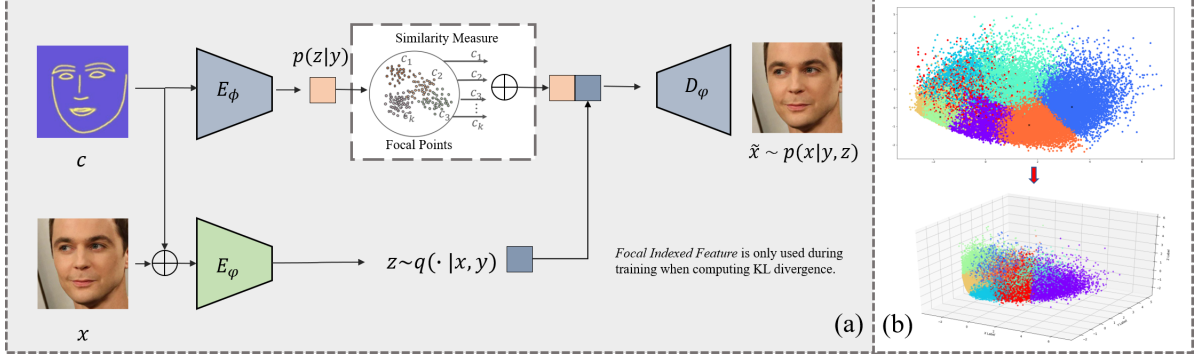
Figure 2: (a) Our framework (b) 2D and 3D projection of 5000 facial structure representations. Each color denotes a cluster, a more intuitive illustration of our approach is to map each cluster to a Gaussian prior, extending the capacity of C-VAE constrained by single prior.

$E_\mu$ and $D_\varphi$ are further incorporated to potentially ensure enough structure information obtained from the prior (*i.e.*, the decoder $D_\varphi$). Thus, $z$ is encouraged to encode more information about appearance rather than structure, otherwise, a penalty of the likelihood $p(x|y, z)$ will be incurred for large reconstruction error.

## 3.2. Additive Memory Encoding

While the explicit geometry-guided disentanglement helps preserving facial structure and structuring latent appearance representation to be invariant towards geomorphing, it is still hard to meet the requirement of fine-grained face manipulation. As shown in Fig. 8, synthesized results on profile with extreme expression or under in-the-wild environment could easily collapse into a "mean-face" phenomenon. The C-VAE formulation is restricted to draw the prior $p(z|y)$ from a simple structure, typically a zero-mean unit variance Gaussian. However, factors like ages, complexions, luminance, and poses constitute a complex distribution, which breaks the common assumption. Consequently, some peculiar facial features might become "outliers" of the distribution during training.

To this end, a better choice of the prior structure is to be explored for a decent appearance representation. Intuitively, drawing the complex prior from a combination of simple distributions will increase the diversity of latent code $z$ and meanwhile enable computing the closed form.

We thus encourage the appearance representation $z$ to have a multi-modal structure composed of $K$ clusters, each corresponding to a semantic feature. In practice, we construct the memory bank through K-means clustering on all boundaries in training set. Every cluster center is called a *focal point*, which commonly refers to a distinctive characteristic such as laughing or side face. In this way, each boundary map used in training would have a k-dimension *focal indexed feature*: $w(b) = (w_1(b), w_2(b), \cdots, w_k(b))$, denoting its similarity measurement with each focal point. The focal indexed feature is used in training to boost the

diversity and plentiful appearance in the latent representation. Since the external memory contains large geometry variation, the general idea is to build up latent representation capacity leveraging explicit and concise spatial semantic guidance provided by focal points. By modeling $p(z|y)$ as a Gaussian Mixture, for each cluster $k$ with weight $w_k$, mean $\mu_k$ and standard deviation $\sigma_k$, we have:

$$p(z|y) = \sum_{k=1}^{K} w_k \mathcal{N}(z|\mu_k, \sigma_k^2 I). \tag{3}$$

However, as it is not directly tractable to optimize Equation 2 with GMM prior, the KL divergence needs to be approximated in training, sampling $z$ from one of the clusters according to their probability. Hence this operation fails to model targets that contain more than one facial geometry characters, *e.g.*, laughing and side face simultaneously. Consequently, we introduce an Additive Focal prior to our framework with the formulation as:

$$p(z|y) = N(z|\sum_{k=1}^{K} w_k \mu_k, \sigma^2 I) \tag{4}$$

where $\sigma^2 I$ is a co-variance matrix with $\sigma^2 = \sum_{k=1}^{K} w_k^2 \sigma_k^2$. Behind the formula, we assume the face image contains multiple structural characteristics with weights $w_k$. The value corresponds to similarity measurement with each focal point, which is defined by a normalized cosine distance. The means $\mu$ of clusters are randomly initialized on the unit ball. The KL term could be computed over $q(z|x, y) = N(z|\mu(x, y), \sigma^2(x, y)I)$, which can be derived to be:

$$D_{KL} = \log(\frac{\sigma}{\sigma_\phi}) + \frac{1}{2\sigma^2} E_{q_\phi}[(z - \sum_{k=1}^{K} w_k \mu_k)^2] - \frac{1}{2}$$
$$= \log(\frac{\sigma}{\sigma_\phi}) + \frac{\sigma_\phi^2 + (\mu_p hi - \sum_{k=1}^{K} w_k \mu_k)^2}{2\sigma^2} - \frac{1}{2} \tag{5}$$

By combining the above KL term into Equation 2, we obtain

the final loss function to train the generator.

## 3.3. Quality-Aware Synthesis

Given a source image $x$ and target boundary $\tilde{b}$, we are able to manipulate the facial expression and head poses of $x$ through the proposed AF-VAE. Following [6, 11], we use feature matching loss $L_{rec} = \|x - \hat{x}\|_1 + \sum_l \lambda_l \|\psi_l(x) - \psi_l(\hat{x})\|$ as reconstruction loss to overcome the blurry results of $L_1/L_2$ losses. However, as can be seen from Table 3, the results are still far from being perceptually appealing.

There are many factors that may cause artifacts in synthesized images and result in perceptual quality distortion, such as losses, unstable training process, and network architectures. Intuitively, one can ease the problem by incorporating auxiliary objective functions [13, 30], or designing sophisticated attention mechanism to capture better global structure [32]. However, we turn to two light-weight designs in network structure through observations and HVS basis for the first time.

Instead of using deconvolution operation to perform upsampling in decoder, we use sub-pixel convolution [38] in every upsamping layer. As shown in Fig. 3, the 'checkerboard' artifacts are reduced apparently. However, it is interesting that the reconstruction results generated from AF-VAE model with/without subpixel yield similar distributions over their histograms, and the differences of their entropy are small (7.6405 without subpixel, 7.6794 with subpixel, and 7.8267 for the source image). It means that subpixel does not filter the artifacts substantially, it scatters them into parts of the image instead. However, HVS treats artifact signals unevenly according to areas of images. Artifacts lie in low frequency or edge areas are emphasized and ones lie in high frequency areas are tend to be masked, as demonstrated in the task of perceptual metrics learning [31, 24, 25, 56]. Thus, we can draw a reasonable explanation for subpixel convolution on the improvements of perceptual measurements.

Another technique is that we utilize Weight Normalization(WN) [36] for model training. It is habitual in GAN-based methods to replace BatchNorm(BN) [19] to WN or its variant for stabilizing the training of the *discriminator*. Instead, we tailor this strategy to both *encoder* and *decoder* of AF-VAE. Models trained with weight normalization have faster convergence and lower reconstruction loss than the ones without WN over training iterations. Detailed analysis and loss curve plot are provided in the appendix. This phenomenon is similar to GANs. However, as can be seen from Fig. 8, WN helps increase the diversity of generated images in complex datasets, making the synthesized results more reasonable for human observation, and thereby improves the perceptual quality. Table 3 shows that the perceptual quality of synthesized results with WN is much better than models with BN and without BN, bringing $57.7\%$
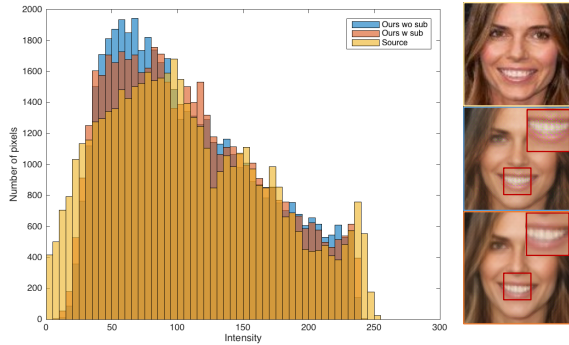


Figure 3: Histograms w/wo pixel-shuffle (Better zoom in).

and $43.6\%$ boosts to IS and FID respectively.

## 4. Experiments

Our framework provides a flexible way to manipulate an input face image into an arbitrary expression and pose under the control of boundary maps. In this section, we show qualitative and quantitative comparisons with state-of-the-art approaches in Sec 4.1. Then, we perform self-evaluation to analyze the key components of our model in Sec 4.2. Finally, we discuss the limitation of our approach in Sec 4.3. All experiments are conducted using the model output from unseen images during the training phase.

**Implementation Details.** Before training, all images are aligned and cropped to $256 \times 256$ resolution. Facial landmarks of each image are obtained using an open-source pre-trained model. Then, landmarks are interpolated to get the facial boundary map. For the clustering, we choose $k = 8$ as the number of clusters and K-means as the clustering algorithm to get clusters of boundaries. For more details, please refer to the appendix.

**Datasets.** We mainly conduct experiments on RaFD [23], MultiPIE [12], and CelebA [26] datasets that cover both indoor and in-the-wild setting. A 3D synthesized face dataset is also introduced to further evaluate the performance of the proposed method on facial texture detail, *e.g.*, illumination, complexion, and wrinkle. We use **hand** and **cat** as non-human datasets. The landmarks of hand and cat dataset are obtained through pre-trained hand detector and human annotation, respectively.

For each datasets, $90\%$ identities are used for training. and the left $10\%$ are fed into the model for testing. Additional quantitative and qualitative results on other datasets are shown in the appendix due to space limits.

**Baselines.** We compare our model with three state-of-the-art GAN-based algorithms: StarGAN [9], GANimation [32], and pix2pixHD [47]. For a fair comparison, we train these models using the implementations provided by the authors, and adopt publicly available pre-trained models to obtain conditional input, like Action Units (AUs) or landmarks for training and testing the models.
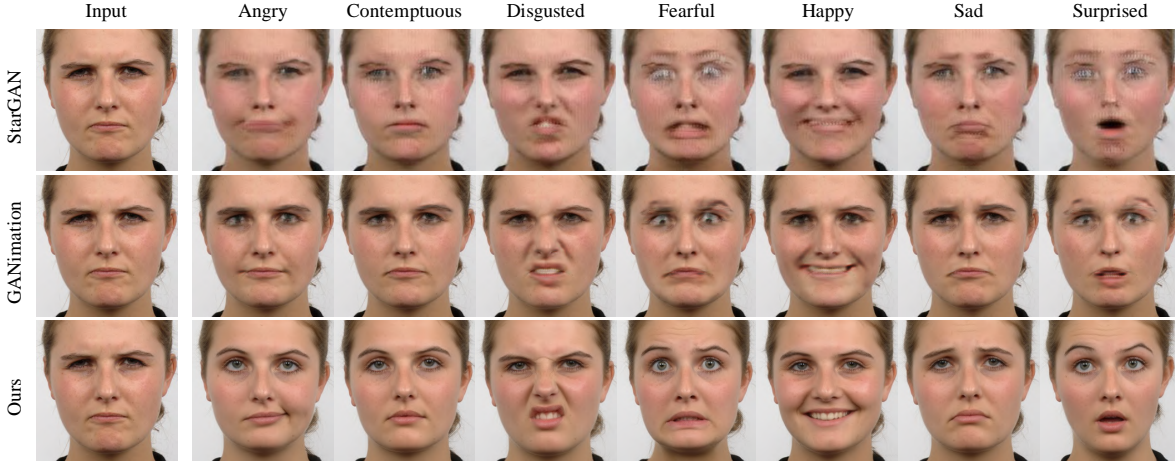
Figure 4: Comparison with three state-of-the-art algorithms.(Zoom in for better details. Left three pictures are input faces, the right three lines are generated results of the three algorithms respectively. StarGAN [9] and GANimation [32] are best current approaches under large poses, still show certain level of blur. Better zoom in to see the detail.
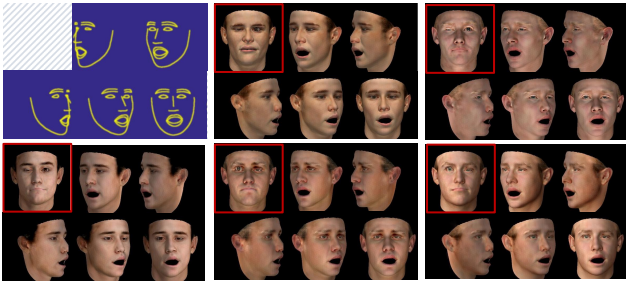


Figure 5: Face rotation results on 3D synthetic face dataset.There are 6 blocks, the upper left block represents 5 target boundary maps. For each blocks, the upper left face with red box is the input image and the rest 5 are synthesized results corresponding to 5 boundary maps in the first block. Each source has different texture and lighting. Better zoom in to see the detail.
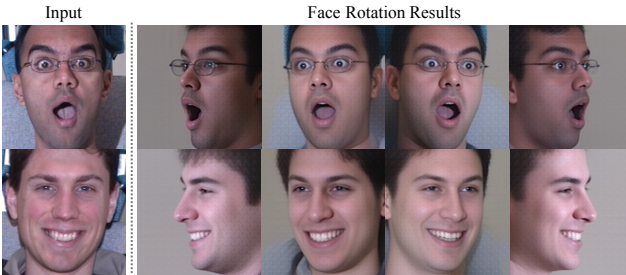


Figure 6: Face rotation results on MultiPIE [12] dataset.

**Performance metrics.** For quantitative comparison, we evaluate three aspects, *i.e*. the realism, the perceptual quality and the diversity, for the synthesis results. For the measurement of perceptual quality and diversity, we use **Fréchet Inception Distance** (FID, lower value indicates better quality) [15] and **Inception Score** (IS, higher value indicates better quality) [35] metrics. For realism, we use Amazon Mechanical Turk (AMT) to compare the perceived visual fidelity of our method against existing approaches.

We report **TS** (TrueSkill) [14] and **FR** (Fool Rate, the estimated probability that generated images succeed in fooling the user) using data gathered from 25 participants per algorithm. Each participant is asked to complete 50 trials.

## 4.1. Comparison with Existing Work
### 4.1.1 Qualitative Comparisons

**Facial Expression Editing.** We conducted face manipulation in comparison with StarGAN, GANimation, and pix2pix on Rafd with 7 typical expressions. As shown in Fig. 4, previous leading methods are fragile when dealing with an exaggerated expression such as "disgusted" and "fearful" with $256 \times 256$ resolution. On the contrary, our method can get rid of blurry artifacts as well as maintaining facial details, owing to finely disentangled latent space and quality-refinement schema. It is worth noting that our results are far better than all of the GAN-based baselines especially in the detailed texture of mouth and eyes, which bring much higher perceptual quality. Quantitative evaluation in Sec. 4.1.2 further demonstrates this observation.

**Face Rotation.** We validate the ability of our model to face rotation task with an arbitrary oriented pose. Note that there is no requirement of any paired training samples and external supervision in our setting, which is different from the state-of-the-art methods such as [18, 16]. First, we conducted qualitative experiments on 3D synthetic face dataset to evaluate the detailed light/texture-preserving nature of our model. The driving landmark maps of different directions are obtained by manipulating the 3D landmark map to new expression and re-projecting it into different 2D views. As shown in Fig. 5, factors including complexion, texture, and lighting can be well preserved. This phenomenon testify the effectiveness of disentangling mechanism on our model. Next, we evaluate the effectiveness on real dataset.

| Model | FID | IS |
|---|---|---|
| Real Data | 0.000 | 1.383 |
| pix2pixHD [47] | 75.376 | 0.875 |
| StarGAN [9] | 56.937 | 1.036 |
| GANimation[32] | 34.360 | 1.112 |
| **Ours** | **25.069** | **1.237** |

Table 1: Quantitative comparison with state-of-the-art on RaFD dataset with FID and IS metrics.

As shown in Fig. 6, even under 90°, our model could still generate high fidelity and photo-realistic results in a simple weakly-supervised fashion.
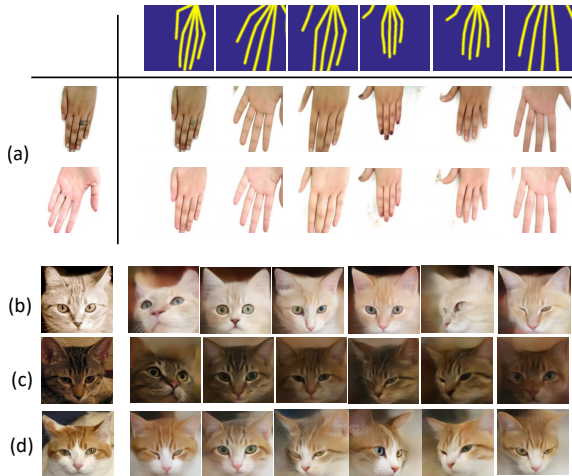


Figure 7: Experiments on Hands and Cats datasets. (a) shows hands manipulation results. On the left are source hands. On the right, the first row presents target hand skeleton and the next two rows represent generated samples respectively. (b)(c)(d) are manipulated results of three cats. Each row presents a source cat image and 6 manipulated results.

### 4.1.2 Quantitative Comparisons

First, we evaluate the perceptual quality and diversity of the generated images. As shown in Table 1, our method outperforms the current state-of-the-art methods by a large margin on both measurements. Our FID and IS is 1.3 and 1.1 times better than the previous leading method, respectively.

Then, we use Amazon Mechanical Turk (AMT) to compare the perceived visual reality of our method against existing approaches. Table 2 reports results of AMT perceptual realism task. We find that our method can fool participants significantly better than other methods. Also, regarding TrueSkill, our model is more likely to gain users' preference, which further supports that our method surpasses those baselines by a large gap in terms of generated reality.

### 4.2. Ablation Study

We conduct ablation study to analyze the contribution of individual components in the proposed method.

| Model | Fool Rate (%) | TrueSkill |
|---|---|---|
| StarGAN [9] | $3\% \pm 0.4\%$ | $18.1 \pm 0.9$ |
| pix2pixHD [47] | $4.8\% \pm 0.9\%$ | N/A |
| GANimation [32] | $7.0\% \pm 1.2\%$ | $24.4 \pm 0.8$ |
| **Ours** | **$36.4\%$** $\pm 2.8\%$ | **$32.6$** $\pm 0.9$ |

Table 2: Evaluation of User Study, compare with state-of-the-arts

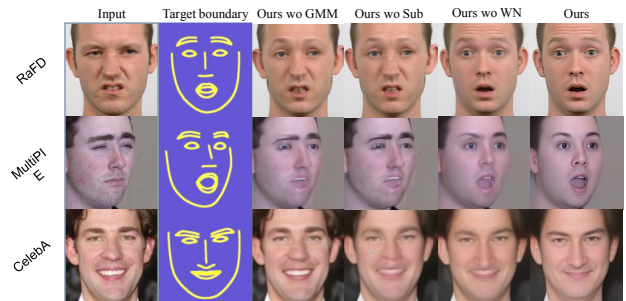| Model | FID | IS |
|---|---|---|
| Real Data | 0.000 | 2.662 |
| Ours w/o KL | 56.275 | 1.863 |
| Ours w/o GMM | 62.657 | 1.792 |
| Ours w/o PS | 71.309 | 1.065 |
| Ours w/o WN | 65.309 | 1.365 |
| **Ours** | **36.820** | **2.152** |

Table 3: Ablation study on CelebA.



Figure 8: Ablation study on RafD, MultiPIE and CelebA datasets.

**Qualitative ablation.** As shown in Fig. 8, without GMM and KL, target with large spatial movements such as opening the mouth and turning faces around could not be generated. Without pixel shuffle, there are artifacts on generated faces. Without WN, the visual quality will be constrained.

**Quantitative ablation.** We evaluate each variant based on the quality of their generated samples with FID and IS metrics. As shown in Table 3, the introduced GMM, pixel shuffle and weight normalization bring great improvement in terms of FID and IS score. Inferred from the above observations, each component has a different role in our method. Removing any of them leads to a performance drop.

**Interpolation Results.** In order to validate that the feature distribution our model learned is dense and distinct, both the appearance and structure in the generated images should change continuously with the latent vector respectively. Our model demonstrates photo-realistic results in Fig.9 via interpolation between disparate samples, suggesting that it has great generalization and robustness, instead of simply memorizing training data.

**Identity Preserving Problem.** Identity preserving is a long standing challenge on face generation domain. Previous leading methods [2, 37] usually address this problem by extending the network with an identity classifier to constrain the variety of synthesized face. To study the influence of different factors to facial identity on our framework, we conduct three experiment settings on RafD, 3D synthetic

Figure 9: Interpolation results on CelebA. Upper left and lower right are two real images, each row and line represents linear Interpolation on latent appearance and structure vector, respectively.

face and EmotionNet datasets, as shown in Fig. 10. We find that when modifying the landmarks from source image (Fig. 10(a)), or using a boundary with similar facial contour from other person (Fig. 10(b)) as the conditional boundary map, the model could synthesize faces with both well preserved identity and high fidelity. Conversely, when the facial contour of conditional boundary map is significantly different with source image(Fig. 10(c)), the identity of synthesized result tend to be diverse. These observation also indicates that facial identity information is mainly encoded in structure space, as demonstrated in [5, 10].

## 4.3. Limitations and Failure Cases

We show four categories of failure cases in Fig. 11, all of them are representative cases that challenge the limit of our model. Specifically, the first case (the top left one) is related to rare data. As cakes are seldom seen in our training samples, the model may not be able to maintain sample's semantics and tend to be blurry. A similar problem occurs in the top right picture where some source parts in source image are occluded. Our model tends to be confused when source image are occluded and may simply move the occlusions from the source due to strong structure coherence in model design. Another challenging case that our model couldn't handle well is special style. As shown in Fig. 11 lower left, characteristics from the source style are likely to be lost when referred to a boundary from human face. As the estimated landmarks may strongly correlate to some attributes such as gender and head pose, simply imposed with
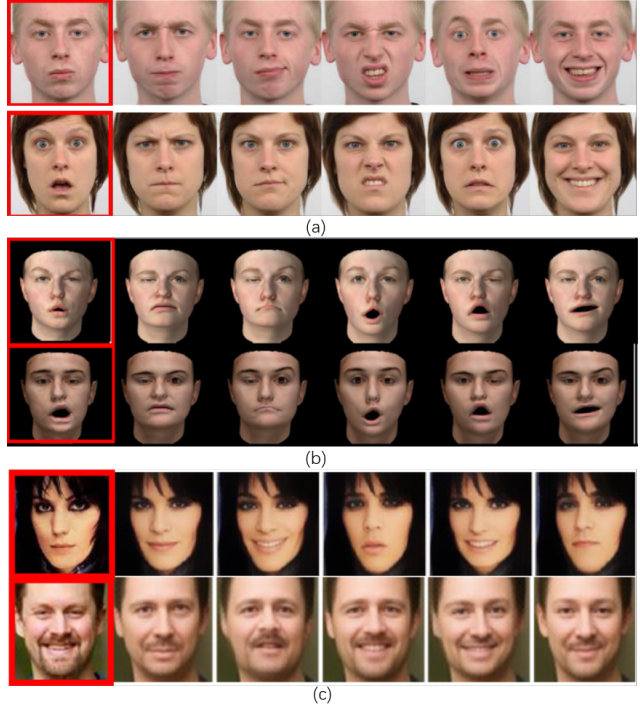


Figure 10: Identity preserving evaluation on three experiment settings.

such a given boundary, the result image are likely to simply morph those attributes with our target appearance as shown in lower right part in Fig. 11. By simply morphing the structure, the target result are unnatural and losing identity.
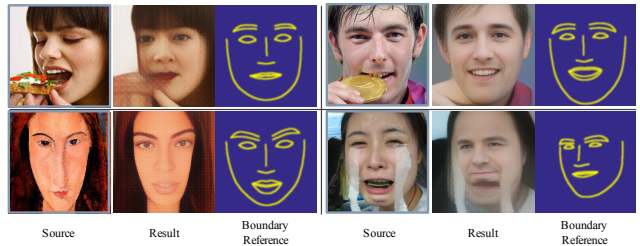


Figure 11: Failure Cases. All four failure cases are selected from CelebA and EmotionNet dataset. We represent the source image at left, followed by manipulation result and its boundary reference.

## 5. Conclusion

In this paper, we propose an additive focal variational auto-encoder (AF-VAE) framework for face manipulation which is capable of modeling the complex interaction between facial structure and appearance. A light-weight architecture designed on HVS basis empowers better synthetic results. It advances current works in face synthesis on both generation quality and adaptation to extreme manipulation settings, as well as its simple structure and stable training procedure. We hope our work could illuminate more trails in this direction.

# References

[1] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua. CVAE-GAN: fine-grained image generation through asymmetric training. In *ICCV*, 2017. 3

[2] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua. Towards open-set identity preserving face synthesis. In *CVPR*, 2018. 2, 7

[3] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999. 2

[4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 3

[5] Chen Cao, Yanlin Weng, Shun Zhou, Yiying Tong, and Kun Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Trans. Vis. Comput. Graph.*, 20(3):413–425, 2014. 8

[6] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. In *ICCV*, 2017. 5

[7] Ying-Cong Chen, Huaijia Lin, Michelle Shu, Ruiyu Li, Xin Tao, Xiaoyong Shen, Yangang Ye, and Jiaya Jia. Faceletbank for fast portrait manipulation. In *CVPR*, 2018. 2

[8] Zeyuan Chen, Shaoliang Nie, Tianfu Wu, and Christopher G. Healey. High resolution face completion with multiple controllable attributes via fully end-to-end progressive generative adversarial networks. *arXiv preprint arXiv:1801.07632*, 2018. 3

[9] 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. 2, 5, 6, 7

[10] Jiankang Deng, Shiyang Cheng, Niannan Xue, Yuxiang Zhou, and Stefanos Zafeiriou. Uv-gan: Adversarial facial uv map completion for pose-invariant face recognition. In *CVPR*, 2018. 8

[11] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A variational u-net for conditional appearance and shape generation. In *CVPR*, 2018. 5

[12] Ralph Gross, Iain Matthews, Jeffrey Cohn, Takeo Kanade, and Simon Baker. Multi-pie. *Image and Vision Computing*, 2010. 5, 6

[13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *NIPS*, 2017. 5

[14] Ralf Herbrich, Tom Minka, and Thore Graepel. Trueskill: a bayesian skill rating system. In *NIPS*, 2007. 6

[15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, 2017. 2, 6

[16] Yibo Hu, Xiang Wu, Bing Yu, Ran He, and Zhenan Sun. Pose-guided photorealistic face rotation. In *CVPR*, 2018. 2, 6

[17] Huaibo Huang, Zhihang Li, Ran He, Zhenan Sun, and Tieniu Tan. Introvae: Introspective variational autoencoders for photographic image synthesis. *arXiv preprint arXiv:1807.06358*, 2018. 3

[18] Rui Huang, Shu Zhang, Tianyu Li, and Ran He. Beyond face rotation: Global and local perception gan for photorealistic and identity preserving frontal view synthesis. In *ICCV*, 2017. 6

[19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 5

[20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 3

[21] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2, 3

[22] Jean Kossaifi, Linh Tran, Yannis Panagakis, and Maja Pantic. Gagan: Geometry-aware generative adversarial networks. In *CVPR*, 2018. 2

[23] Oliver Langner, Ron Dotsch, Gijsbert Bijlstra, Daniel HJ Wigboldus, Skyler T Hawk, and AD Van Knippenberg. Presentation and validation of the radboud faces database. *Cognition and emotion*, 2010. 5

[24] Kwan-Yee Lin and Guanxiang Wang. Hallucinated-iqa: No-reference image quality assessment via adversarial learning. In *CVPR*, 2018. 5

[25] Xialei Liu, Joost van de Weijer, and Andrew D. Bagdanov. Rankiqa: Learning from rankings for no-reference image quality assessment. In *ICCV*, 2017. 5

[26] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, 2015. 2, 5

[27] Yongyi Lu, Yu-Wing Tai, and Chi-Keung Tang. Attribute-guided face generation using conditional cyclegan. In *ECCV*, 2018. 2

[28] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. In *NIPS*, 2017. 2

[29] Liqian Ma, Qianru Sun, Stamatios Georgoulis, Luc Van Gool, Bernt Schiele, and Mario Fritz. Disentangled person image generation. In *CVPR*, 2018. 2

[30] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *ICCV*, 2017. 5

[31] Da Pan, Ping Shi, Ming Hou, Zefeng Ying, Sizhe Fu, and Yuan Zhang. Blind predicting similar quality map for image quality assessment. In *CVPR*, 2018. 5

[32] 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. 2, 5, 6, 7

[33] Shengju Qian, Wayne Wu, Yangxiaokang Liu, Beier Zhu, and Fumin Shen. Extending the capacity of cvae for face synthesis and modeling. In *NeurIPS Workshops*, 2018. 3

[34] Fengchun Qiao, Naiming Yao, Zirui Jiao, Zhihao Li, Hui Chen, and Hongan Wang. Geometry-contrastive generative adversarial network for facial expression synthesis. *arXiv preprint arXiv:1802.01822*, 2018. 2

[35] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *NIPS*, 2016. 6

[36] Tim Salimans and Diederik P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *NIPS*, 2016. 5

[37] Yujun Shen, Ping Luo, Junjie Yan, Xiaogang Wang, and Xiaoou Tang. Faceid-gan: Learning a symmetry three-player gan for identity-preserving face synthesis. In *CVPR*, 2018. 2, 7

[38] 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. 5

[39] Zhixin Shu, Mihir Sahasrabudhe, Riza Alp Güler, Dimitris Samaras, Nikos Paragios, and Iasonas Kokkinos. Deforming autoencoders: Unsupervised disentangling of shape and appearance. In *ECCV*, 2018. 3

[40] Zhixin Shu, Ersin Yumer, Sunil Hadap, Kalyan Sunkavalli, Eli Shechtman, and Dimitris Samaras. Neural face editing with intrinsic image disentangling. In *CVPR*, 2017. 3

[41] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *NIPS*, 2015. 2

[42] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *CVPR*, 2016. 2

[43] Luan Tran and Xiaoming Liu. Nonlinear 3d face morphable model. In *CVPR*, 2018. 2

[44] Luan Tran and Xiaoming Liu. On learning 3d face morphable model from in-the-wild images. *arXiv preprint arXiv:1808.09560*, 2018. 2

[45] Liwei Wang, Alexander Schwing, and Svetlana Lazebnik. Diverse and accurate image description using a variational auto-encoder with an additive gaussian encoding space. In *NIPS*, 2017. 3

[46] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *NIPS*, 2018. 2

[47] 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. 3, 5, 7

[48] Wei Wang, Xavier Alameda-Pineda, Dan Xu, Pascal Fua, Elisa Ricci, and Nicu Sebe. Every smile is unique: Landmark-guided diverse smile generation. In *CVPR*, 2018. 2

[49] Olivia Wiles, A Koepke, and Andrew Zisserman. Self-supervised learning of a facial attribute embedding from video. In *BMVC*, 2018. 2

[50] Olivia Wiles, A Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *ECCV*, 2018. 2

[51] Wayne Wu, Kaidi Cao, Cheng Li, Chen Qian, and Chen Change Loy. Transgaga: Geometry-aware unsupervised image-to-image translation. In *CVPR*, 2019. 2

[52] Wayne Wu, Chen Qian, Shuo Yang, Quan Wang, Yici Cai, and Qiang Zhou. Look at boundary: A boundary-aware face alignment algorithm. In *CVPR*, 2018. 3

[53] Wayne Wu, Yunxuan Zhang, Cheng Li, Chen Qian, and Chen Change Loy. Reenactgan: Learning to reenact faces via boundary transfer. In *ECCV*, 2018. 2

[54] Ceyuan Yang, Zhe Wang, Xinge Zhu, Chen Huang, Jianping Shi, and Dahua Lin. Pose guided human video generation. In *ECCV*, 2018. 2

[55] Xi Yin, Xiang Yu, Kihyuk Sohn, Xiaoming Liu, and Manmohan Chandraker. Towards large-pose face frontalization in the wild. In *ICCV*, 2017. 2

[56] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5