

# MIMAFace: Face Animation via Motion-Identity Modulated Appearance Feature Learning

Yue Han<sup>1</sup>, Junwei Zhu<sup>2</sup>, Yuxiang Feng<sup>1</sup>, Xiaozhong Ji<sup>2</sup>, Keke He<sup>2</sup>,  
Xiangtai Li<sup>3</sup>, Zhucun Xue<sup>1</sup>, Yong Liu<sup>†1</sup>

<sup>1</sup>Zhejiang University <sup>2</sup>Tencent Youtu Lab <sup>3</sup>Nanyang Technological University  
{12432015, fengyx, 12432038}@zju.edu.cn, yongliu@iipc.zju.edu.cn

<https://mimaface2024.github.io/mimaface.github.io>

## Abstract

Current diffusion-based face animation methods generally adopt a ReferenceNet (a copy of U-Net) and a large amount of curated self-acquired data to learn appearance features, as robust appearance features are vital for ensuring temporal stability. However, when trained on public datasets, the results often exhibit a noticeable performance gap in image quality and temporal consistency. To address this issue, we meticulously examine the essential appearance features in the facial animation tasks, which include motion-agnostic (e.g., clothing, background) and motion-related (e.g., facial details) texture components, along with high-level discriminative identity features. Drawing from this analysis, we introduce a Motion-Identity Modulated Appearance Learning Module (MIA) that modulates CLIP features at both motion and identity levels. Additionally, to tackle the semantic/ color discontinuities between clips, we design an Inter-clip Affinity Learning Module (ICA) to model temporal relationships across clips. Our method achieves precise facial motion control (i.e., expressions and gaze), faithful identity preservation, and generates animation videos that maintain both intra/inter-clip temporal consistency. Moreover, it easily adapts to various modalities of driving sources. Extensive experiments demonstrate the superiority of our method.

## 1. Introduction

Face Animation aims to generate a realistic talking head by animating a source face using the motion information of another person, i.e., pose, expression, and gaze [23]. It has diverse applications in virtual character creation for game production and video editing. Previous GAN-based methods [2, 14, 15, 22, 27, 29, 30, 33, 43, 47, 49, 51, 53] most delivers results at resolution 256<sup>2</sup> and only support median

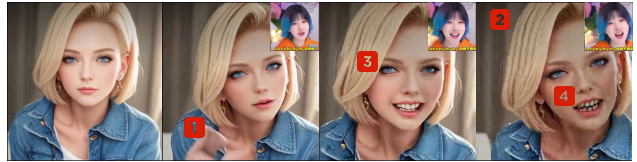


Figure 1. Typical failure cases for current diffusion-based face animation methods: (1)/(2) semantic/ color discontinuity across clips, (3) stiff expression, (4) quality degradation

pose variation, i.e., less than 30 degrees), due to the lack of high-resolution and large pose dataset. Recent attempts leverage the powerful generation capability of pre-trained latent diffusion models to address these challenges. However, the high variance of noise in diffusion presents a new challenge in generating smooth videos.

The human-body animation methods [17, 46] identify that CLIP [26] fails to provide adequate appearance features, resulting in video flickering. To address this issue, ReferenceNet (a UNet copy) is proposed to supply multi-scale similar appearance features. Subsequent works in face animation adopt similar frameworks, integrating ReferenceNet with motion modules proposed in Animatediff [11] to ensure temporal stability and achieve remarkably favorable results. However, the quality of the results within this paradigm [4, 17, 34, 38, 46] heavily depends on the quality and volume of the training data. When trained on publicly available datasets [4, 38], the generated results often exhibit a noticeable gap in quality, as shown in Fig. 1 (from released examples in AniPortrait). A question arises: *what appearance features are essential for producing a temporally stable, high-quality face animation video?* We argue that the necessary appearance features include motion-independent texture features (i.e., clothing, background), motion-related texture features (i.e., facial details), and high-level discriminative features (i.e., identity). Based on this analysis, we propose the Motion-Identity Modulated Appearance Learning Module (MIA), which modulates CLIP features at both motion and identity levels. For

\*<sup>†</sup>denotes corresponding author.

motion modulation, we use 3DMM coefficients to modulate appearance features via cross-attention. This facilitates the generation of subtle facial textures, *e.g.*, wrinkles and muscle contractions, resulting from expressions. For identity modulation, we introduce an identity contrastive loss to compensate for high-level discriminative features lost. This is because we find that the optimization objective, *i.e.*, the denoise loss, encourages CLIP to focus more on the learning of low-level generative features while neglecting the learning of discriminative features. This necessitates that the model learns to preserve identity through training on data that encompasses a wider range of identities. The joint modulation of motion and identity makes appearance features learning easier, allowing us to achieve high-fidelity and temporally stable video results even when trained solely on public datasets.

Although adequate appearance features ensure intra-clip temporal consistency with temporal attention, the lack of modeling inter-clip relationships causes semantic/ color discontinuities between clips. Current solutions include fixing initial noise, frame interpolation, or temporal co-denoising, but noticeable discontinuities still exist. To address this issue, we introduce the Inter-clip Affinity Learning Module (ICA), which conditions preceding frames with augmentation strategies to bridge the gap between ground truth frames during training and generated frames during inference.

Our contributions can be summarized as follows:

- We identify potential issues in current diffusion-based face animation methods and carefully examine the essential appearance features. Based on the analysis, we introduce a Motion-Identity Modulated Appearance Learning Module (MIA) that modulates CLIP features at both motion and identity levels, making the appearance features learning more effective.
- We introduce the Inter-clip Affinity Learning Module (ICA), which models temporal relationships between clips to address the issue of semantic/color discontinuities.
- Our method achieves precise facial motion control (*i.e.* expressions and gaze), faithful identity preservation, and generates animation videos that maintain both intra/inter-clip temporal consistency. Moreover, it easily adapts to various modalities of driving sources. Extensive experiments demonstrate the superiority of our method.

## 2. Related Works

### 2.1. GAN-based Face Animation

Methods [1–3, 8, 22, 30, 43, 44, 47, 51, 53] can be broadly divided into warping-based and 3DMM-based methods. *Warping-based methods* [15, 29, 30, 58] typically extract

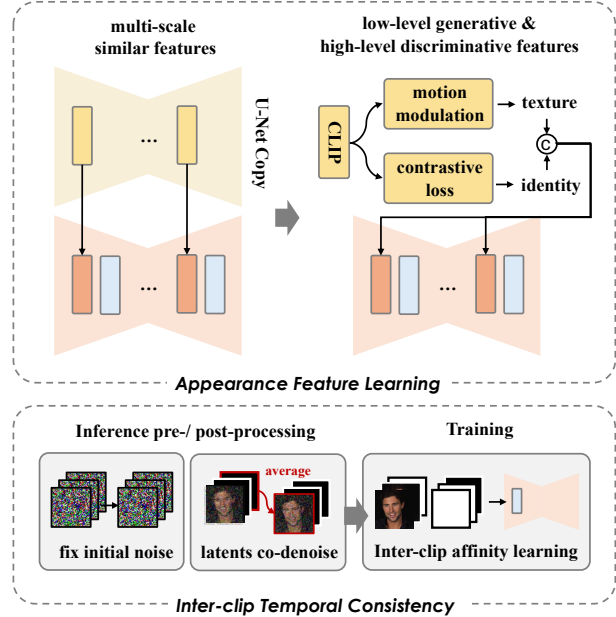


Figure 2. We compare our method to previous diffusion-based face animation methods in terms of appearance feature learning and inter-clip temporal consistency.

landmarks or region pairs to estimate flow fields and perform warping on the source appearance feature maps to transfer motions. Limited by the accuracy of the predicted flow field, these methods tend to produce blurry and distorted results when dealing with large motion variations.

*3DMM-based methods* use facial reconstruction coefficients or the render image from 3DMM as motion intermediate representation. Due to the inherent decoupling properties of coefficients, 3DMM-based methods, *e.g.*, PIRenderer [27], can freely control expressions and poses. Although 3DMM provides accurate structural references for facial regions, it lacks references for hair, teeth, and eye movement. Additionally, its coarse facial textures result in suboptimal generated outcomes. StyleHEAT [49] tackles this challenge by leveraging the powerful generation capabilities of StyleGAN2 to produce detailed textures and high-resolution frontal portraits. However, its generalization and large pose generation abilities are limited due to the constraint of the training dataset.

### 2.2. Diffusion-based Face Animation

To improve sample quality and generalization capability, diffusion models have gained popularity [12, 17, 24, 46, 50]. FADM [50] combines the previous reenactment models with diffusion refinements, but the base model limits the driving accuracy. The human-body animation methods [17, 46] identify that CLIP fails to provide adequate appearance features, resulting in video flickering. To ad-

dress this issue, ReferenceNet (a UNet copy) is proposed to supply multi-scale similar appearance features. Subsequent works in face animation adopt similar frameworks, integrating ReferenceNet with motion modules proposed in AnimateDiff [11] to ensure temporal stability and achieve remarkably favorable results. However, the quality of the results within this paradigm [4, 17, 34, 38, 46] heavily depends on the quality and volume of the training data. When trained on publicly available datasets [4, 38], the generated results often exhibit a noticeable gap in quality. To address this issue, in this paper, we explore the effective way to learn robust appearance features.

### 2.3. Audio-Driven Face Animation

Previous approaches [6, 7, 9, 18, 25, 28, 32, 36, 42, 45, 48, 52, 54, 59] focus on learning models specific to individual speakers. Sadtalker [56] uses 3DMM as an intermediate representation for subject-agnostic reenactment. Via learning the 3D motion coefficients of the 3DMM model from audio, Sadtalker showcases robust generalization capabilities. However, it still struggles to accommodate significant motion variations and produces blurry results. Recently, several works [34, 38] have followed the 'ReferenceNet with motion module' paradigm and achieved significant advancements. [34] bypasses the need for intermediate 3D models or facial landmarks and achieve astonishing results. However, it heavily relies on the dataset, making it difficult to reproduce. [38] learns to map the audio to a 3D facial mesh and head pose, then projects these two elements into 2D keypoints as the intermediate motion representation. However, it also exhibits typical failure cases similar to methods using the same paradigm. In this paper, our primary goal is to explore how to effectively learn robust appearance features, not to focus on the learning of motion intermediate representations. Following the approach of [38, 49], which can be easily adapted to multi-modal inputs, we support audio-driven face animation by mapping audio to the 3DMM coefficient space using pre-trained audio-to-3DMM-coefficients encoders, e.g., [56].

### 2.4. Temporal Consistency in Face Animation.

Current methods [4, 34, 38–40, 46] tend to focus more on modeling the temporal stability within clips while neglecting the preservation of temporal consistency between clips. Recent image animation [57] methods (typically generating only one clip) demonstrate that temporal attention is sufficient for ensuring intra-clip temporal stability. Therefore, we argue the temporal stability within clips is mainly affected by the lack of adequate appearance features. Current solutions for inter-clip temporal stability include fixing initial noise, frame interpolation, or temporal co-denoising [46], but noticeable discontinuities still persist. Unlike these pre/post-processing approaches during

inference, we explicitly model the temporal relationships between clips during training.

## 3. Methods

Face Animation aims to create a lifelike talking head video by animating a source face  $I_S$  using motion information from another person. This motion information may include pose, expression, and gaze, which can be obtained from either a video sequence  $I_D^{1:N}$  or an audio sequence  $A_D^{1:N}$ . The entire sequence has a length of  $N$  and is divided into several clips with the length of  $F$ .

The current challenge in generating high-quality and smooth animation videos using diffusion-based methods can be summarized as the need for robust appearance features and ensuring temporal consistency within and between clips: **1)** Recent methods commonly use ReferenceNet for extracting appearance features, but this approach relies on learning from large amounts of curated data. We believe this is due to the lack of guidance on effective optimization directions. Thus, we propose MIA that modulates appearance features at both motion and identity levels to address this issue. The motion modulation encourages the model to learn facial texture details related to motion, while the identity modulation focuses on high-level discriminative features beyond low-level texture information guided by denoise loss. **2)** To ensure temporal consistency, the motion module is shown to maintain intra-clip temporal consistency effectively. However, the challenge of ensuring inter-clip consistency remains largely unexplored. Existing methods [38, 46] like frame interpolation, fixed initial noise, and temporal co-denoising have proven insufficient in resolving flickering. In this work, we propose ICA to address this issue by modeling relationships with preceding frames during training. The pipeline of the proposed MI-MAFace is illustrated in Fig. 3

### 3.1. Motion Intermediate Representation

We utilize the motion coefficients of 3DMM, *i.e.*, pose  $\rho$ , expression  $\beta$ , and the rendered image  $I_R$ , as a unified motion intermediate representation (for details on 3DMM, please refer to the supplementary materials). The rendered image  $I_R$  is conditioned by concatenating with the noisy latent  $z_t$ . The coefficients are used to modulate the appearance features to supplement more motion-related facial details. This approach offers several advantages: **1)** Highly disentangled coefficient space allows for convenient integration of multi-modal driving sources, mapping images/audio to 3D coefficients with pre-trained models; **2)** Continuous coefficient space facilitates smooth video generation, in contrast to the spatial jittering characteristics of landmarks; **3)** The rendered image  $I_R$ , similar to landmarks, provides coarse spatial guidance, while the coefficient encodes more information, compensating for the missing de-

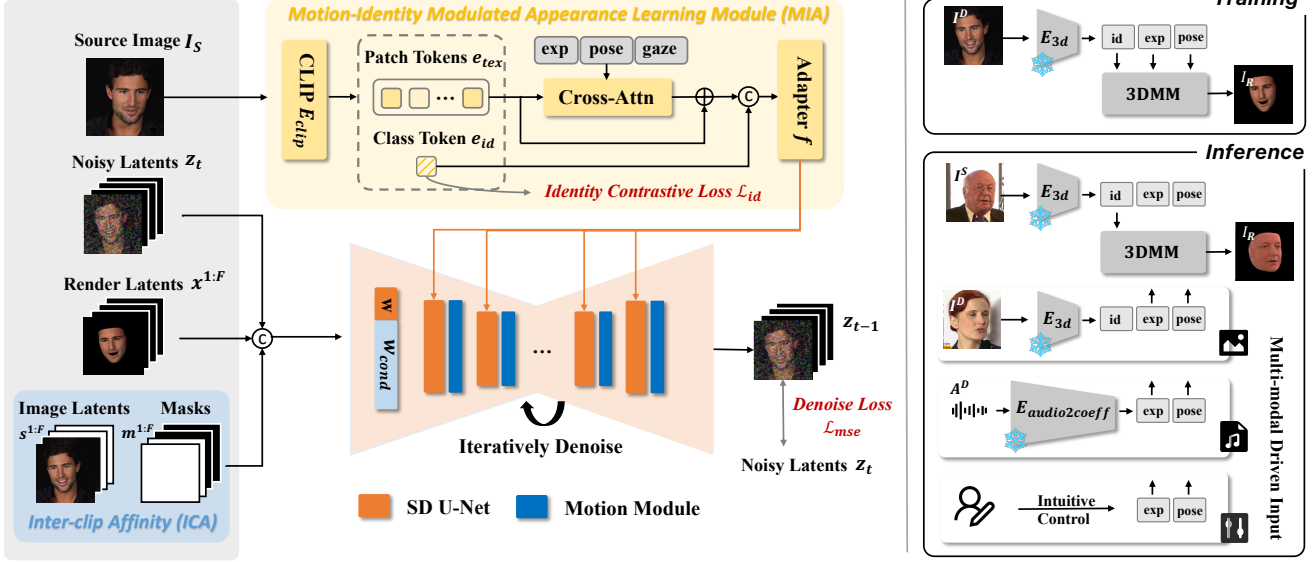


Figure 3. **Pipeline of the proposed MIMAFace**, which consists of: 1) **Motion-Identity Modulated Appearance Learning Module (MIA)** and 2) **Inter-clip Affinity Learning Module (ICA)**. MIA modulates the appearance features at both motion and identity levels. The source image  $I_S$  is passed to CLIP  $E_{clip}$  to obtain patch tokens  $e_{tex}$  and a class token  $e_{id}$ , which capture the texture and identity, respectively.  $e_{tex}$  are then modulated with motion coefficients  $\rho, \beta, g$  via cross attention.  $e_{id}$  is used to calculate the identity contrastive loss  $\mathcal{L}_{id}$ . The modulated  $e_{tex}$  are concatenated with  $e_{id}$  to form the conditioning appearance features. ICA ensures inter-clip temporal consistency by conditioning image latent  $s^{1:F}$  (of ground truth during training and denoised ones during inference) and indicating masks  $m^{1:F}$  with the added condition module  $W_{cond}$ . Additionally, we employ 3DMM coefficients  $\rho, \beta$  and rendered images  $I_R$  as intermediate representations for motion. The 3DMM coefficients can adapt to various modalities of driving inputs, *i.e.*, images, audio, and manual modifications.

tails about subtle facial expressions. Besides 3DMM, we employ an off-the-shelf gaze detector to extract gaze embedding  $g$ .

### 3.2. Motion-Identity Modulated Appearance Learning Module (MIA)

The robust appearance features should include motion-independent texture features (*i.e.*, clothing, background), motion-related texture features (*i.e.*, facial details), and high-level discriminative features (*i.e.*, identity). However, under the guidance of the denoise loss, the appearance encoder tends to learn common textures more easily. This leads to two issues: 1) Although the appearance interacts with the motion intermediate representation through SD cross-attention, it is difficult for landmarks/render to express subtle facial movements. 2) The optimization target for generation enc, we introduce modulation at both the motion and identity levels to improve the effectiveness and robustness of appearance features. To improve the effectiveness and robustness of appearance features, we introduce modulation at both the motion and identity levels. Specifically, we still choose CLIP as the appearance encoder to validate our hypothesis. The source image is passed through the CLIP vi-

sion encoder  $E_{clip}$ , yielding patch tokens  $e_{tex}$  to capture low-level textures, and a class token  $e_{id}$  to represent high-level identity.

**Motion Modulation** To enable the model to generate more vivid and subtle facial movements, we introduce additional 3DMM motion parameters  $\rho, \beta$  and gaze embedding  $g$  to modulate the texture tokens  $e_{tex}$  via simple cross-attention with a residual connection to prevent information loss:

$$e'_{tex} = e_{tex} + \text{Cross-Attn}(e_{tex}, [\rho; \beta; g]),$$

$$\text{Cross-Attn}(e_{tex}, [\rho; \beta; g]) = \text{softmax}\left(\frac{e_{tex} \cdot [\rho; \beta; g]^T}{\sqrt{d_k}}\right) [\rho; \beta; g], \quad (1)$$

where  $d_k$  is the dimensionality of the key. This process can be understood as adding more motion information to refine the coarse spatial structure determined by the render. It can also be seen as aligning the facial appearance features to the same motion as the render, making the cross-attention feature interaction within SD more effective.

**Identity Modulation** To prevent the model from neglecting the learning of high-level features, we introduce a discriminative objective. For faces, identity is the most intuitive high-level feature. Therefore, we incorporate an *Identity Contrastive Loss*  $\mathcal{L}_{id}$  during training, as illustrated in Fig. 4. Specifically, we augment the source image  $I_S$  through



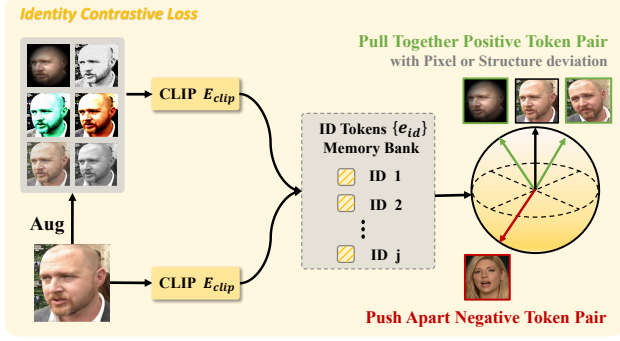


Figure 4. **Illustration of our Identity Contrastive Loss.** We apply photometric data augmentation on the source image and maintain an ID token memory bank to store ID tokens. By pulling together positive token pairs with variances in pixels or structure and pushing apart the negative token pairs, the loss encourages the appearance encoder to capture high-level discriminative features.

photometric transformation, altering the RGB channels by shifting pixel colors to new values. This includes techniques such as grayscaling, color jittering, various filtering methods (like edge enhancement, blurring, and sharpening), lighting perturbation, noise addition, vignetting, and contrast adjustment. Subsequently, the source image  $I_S$  and the augmented image  $I_{aug}$  are each processed through  $E_{clip}$  to obtain a positive ID token pair. During the training process, we maintain an ID token memory bank  $\mathcal{M} = \{e_{id}\}$  to store ID tokens. Consequently, an ID token belonging to the same person but with a different structure will produce another type of positive ID token pair. Explicitly combining the same identity samples across pixel and structural variations enhances the model’s generalization capability and robustness. The Identity Contrastive Loss  $\mathcal{L}_{id}$  is given by:

$$\mathcal{L}_{id} = -\log \left( \frac{\exp(\text{sim}(z_i, z_i^+))}{\exp(\text{sim}(z_i, z_i^+)) + \sum_{j=1}^N \exp(\text{sim}(z_i, z_j^-))} \right), \quad (2)$$

where  $z_i$  represents the token of the identity token,  $z_i^+$  and  $z_j^-$  represent the token of the same and different identity, respectively.  $\text{sim}(z_i, z_j)$  denotes the cosine similarity between two tokens. Finally, the modulated patch tokens  $e'_{tex}$  and the class token  $e_{id}$  are concatenated and go through an adapter  $f$  constructed by a 3-layer transformer to obtain the final appearance condition.

### 3.3. Inter-clip Affinity learning Module (ICA)

Although temporal consistency within video clips is learned and preserved through temporal attention layers, variations in lighting, color, and semantics may still occur between different clips. Previous methods [38, 46] commonly apply noise pre-processing or denoised latent post-processing during inference without explicitly learning temporal correspondence.

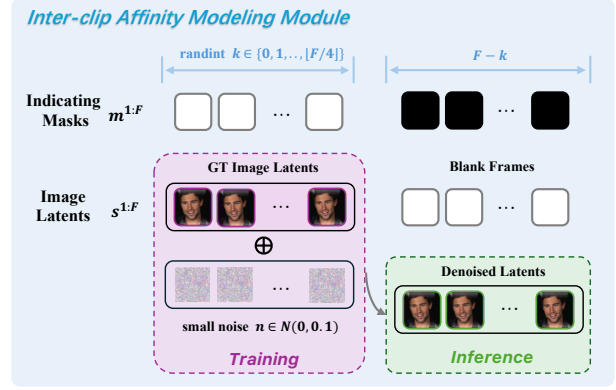


Figure 5. **Illustration of our Inter-clip Affinity Learning Module.** The model learns inter-clip temporal consistency by conditioning the image latent of the preceding frames and using masks to indicate whether reconstruction is required.

Our goal is to enable the model to smoothly generate the next video clip by referring to the last few frames in the previous clip. This problem can be reformulated as follows: When the model is conditioned on the preceding frames, we aim for it to reconstruct these frames, thereby maintaining temporal consistency with the preceding frames. This is because the temporal consistency of subsequent frames within the video clip can be ensured by temporal attention. During training, given a video clip of length  $F$ , the initial  $k$  frames are ground truth image latent  $s_{gt}^{1:k}$ , and the remaining frames  $s^{k+1:F}$  are padded with zeros. Meanwhile, one-channel masks  $m^{1:F}$  are used to indicate whether the model should reconstruct the given image latent correspondingly.

$$s^{1:k} = s_{gt}^{1:k}, \quad s^{k+1:F} = 0, \quad m^{1:k} = 1, \quad m^{k+1:F} = 0, \quad (3)$$

where  $k \sim \text{Uniform}\{0, 1, \dots, \lfloor F/4 \rfloor\}$ ,

This approach differs from using denoised latent codes  $s_d^{1:F}$  during inference, which we find leads to a decrease in video quality. To bridge this domain gap, we introduce a small amount of Gaussian noise to the ground truth image latent during training, simulating the denoised latent:

$$s_{gt, \text{noise}}^{1:k} = s_{gt}^{1:k} + \mathcal{N}(0, \sigma^2), \quad (4)$$

Here,  $\mathcal{N}(0, \sigma^2)$  represents Gaussian noise with mean 0 and standard deviation  $\sigma$  (set to 0.1 here). Inspired by inter-frame affinity in PIA, we condition the image latent  $s^{1:F}$  in the input of the U-Net using a lightweight learnable single-layer convolution  $W_{cond}$  as the condition module, without compromising the original functionality, as shown in Fig. 5.

## 4. Experiment

**Datasets.** We train our model on VoxCeleb2 [21]. We evaluate the VoxCeleb1 test set following the sampling strategy



Figure 6. Same-identity and cross-identity reenactment results on Voxceleb1 test set.

of the PIRenderer [27]. We test on FFHQ [20] to verify the generalization ability.

**Metrics.** We use PSNR and LPIPS [55] to evaluate the reconstruction quality for same-identity reenactment. Expression, pose, and gaze accuracy are assessed by calculating the average Euclidean distance of the corresponding embedding between the generated and driving faces. These three embeddings are derived through the respective estimator. Identity preservation cosine similarity (CSIM) is calculated by [19]. FID is used to evaluate the realism of the generated faces.

**Training Details.** Our training process is divided into two stages. In the first stage, we train the image-driven model. We begin training from the StableDiffusion v1-5 model and OpenAI clip-vit-large-patch14 vision model. Our models are trained for 30k steps on 4 NVIDIA A100 GPUs, with a constant learning rate of  $1e-5$  and a batch size 32. To facilitate classifier-free guidance sampling, we train the model without appearance conditions on 10 of the instances. In the second stage, we train the video-driven model. We freeze the SD U-Net and the appearance encoder and train the condition module  $W_{cond}$ , along with the motion module. Our models are trained for 30k steps on 4 NVIDIA A100 GPUs, with a constant learning rate of  $1e-5$ , a batch size of 8, and a clip sequence of 12.

#### 4.1. Comparison with State-of-the-Art Methods

**Methods.** For image-driven reenactment, we compare our method with GAN-based methods, including FOMM [29], Face-vid2vid [37], PIRenderer [27], TPSM [58], DAM [33] and diffusion-based methods, including FADM [50] and AniPortrait [38]. Except AniPortrait, all of these models are trained on VoxCeleb1. For video-driven reenactment, we compare with MCNet [14].

**Qualitative Results.** In Fig. 6, previous GAN-based methods tend to produce blurry results with noticeable artifacts. While FADM refines the results of FOMM or Face-vid2vid using diffusion, enhancing the generation quality, this approach inherently inherits the motion deviations from these methods. In contrast, our method delivers results of both high fidelity and precise control simultaneously. Compared to previous outcomes that appear relatively smooth, it is worth noting the nuanced skin wrinkles and light-shadow variations caused by expressions in our method, which make the generated results appear much more lifelike.

**Quantitative Results.** In Tab. 1, according to FADM, the data quality of VoxCeleb is relatively subpar, characterized by low resolution and blurred textures. Diffusion-based methods tend to generate fine-detailed images, leading to a mismatch between them. Consequently, we face noticeable disadvantages in pixel-level metrics (PSNR and LPIPS), and distribution-level metrics (FID). However, in



Table 1. Quantitative evaluations among current popular methods on Voxceleb1 test set.

Methods	Same-Identity						Cross-Identity					
	PSNR↑	LPIPS↓	Exp↓	Pose↓	Gaze↓	CSIM↑	FID↓	Exp↓	Pose↓	Gaze↓	CSIM↑	FID↓
FOMM [29]	22.38	0.1405	2.77	0.0261	0.0554	0.8328	26.69	6.28	0.0638	0.0959	0.5642	42.76
PIRenderer [27]	21.00	0.1468	2.94	0.0496	0.0800	0.7997	25.48	5.84	0.0752	0.0977	0.5659	35.99
Face-vid2vid [37]	22.63	<b>0.1243</b>	2.84	0.0283	0.0850	0.8383	<b>25.36</b>	6.68	0.0847	0.1220	0.6328	39.87
TPSM [58]	23.24	0.1442	2.58	0.0224	0.0538	0.8277	33.63	6.10	0.0535	0.0900	0.5836	50.43
DAM [33]	<b>23.37</b>	0.1550	2.81	0.0263	0.0628	0.8333	36.40	6.31	0.0626	0.0967	0.5534	54.13
FADM [50]	22.36	0.1425	2.95	0.0303	0.0879	0.8352	31.70	6.71	0.0821	0.1242	0.6522	42.22
Ours	18.64	0.1907	<b>2.21</b>	<b>0.0195</b>	<b>0.0482</b>	<b>0.8412</b>	60.91	<b>5.03</b>	<b>0.0503</b>	<b>0.0614</b>	<b>0.6778</b>	64.49



Figure 7. Ablation study of different components on the VoxCeleb test set.



Figure 8. Evaluation on FFHQ. Our method exhibits outstanding generalization capabilities. Even when dealing with unseen identities, it maintains consistent identity preservation. Moreover, it ensures intricate facial textures and precise expressions.

terms of semantic-level metrics, including identity similarity (CSIM) and motion accuracy (expression, pose, and gaze), we exhibit clear advantages.

**Handling Unseen Identities.** In Fig. 8, we evaluate on FFHQ to verify the generalization capability of our approach. Whereas previous methods yield blurry results, our approach consistently delivers high-fidelity generation.

## 4.2. Video-driven Reenactment

MCNet is a method designed specifically for videos and fail when driven by a single image or a short video. Thus, we



Figure 9. Video-driven Reenactment Results.

compare our results exclusively with theirs in the context of video generation in Fig. 9. MCNet relies on a memory mechanism to maintain clarity in appearance under continuous large-motion changes in videos. However, if the drive image and source image are initially misaligned, the errors in pose and expression will gradually accumulate, leading to a decline in generated quality and the appearance of noticeable artifacts. In contrast, our approach can produce more accurate motion.

## 4.3. Ablation Studies

We perform ablation studies in Fig. 7 to verify the effectiveness of the Motion-Identity Modulated Appearance Learning Module (MIA).

**Motion Modulation.** If we remove the render image and rely solely on 3DMM coefficients for motion control, it re-

sults in a dramatic decrease in the precision of the generated pose and expression. because the model loses a direct reference to the spatial position, making the learning process more challenging. If we remove the 3DMM coefficients and rely solely on the render image for motion control, it results in a moderate decrease in the precision of the generated pose and expression.

**Identity Modulation.** Without incorporating the identity contrastive loss, the model, during finetuning, is inclined by the denoise loss to capture low-level features, leading to a decline in identity discrimination capability.

## 5. Conclusion

This work proposes a novel MIMAFace framework to address the limitations of current diffusion-based face animation methods. We introduce a Motion-Identity Modulated Appearance Learning Module (MIA) that modulates CLIP features at both motion and identity levels, and an Inter-clip Affinity Learning Module (ICA) to model temporal relationships across clips. Our method ensures precise facial motion control and faithful identity preservation, and generates animation videos with intra/inter-clip temporal consistency. Extensive experiments demonstrate the effectiveness of the proposed method.

**Limitation and Future Works.** Integrating large language models and face-specific SD can enhance performance and application value. However, since our approach can generate realistic images that could be used for facial forgery, regulatory constraints are necessary to mitigate this risk.



## 6. Appendix

### 6.1. Details of 3DMM

We choose motion coefficients and the rendered image obtained by 3D Morphable Models (3DMMs) as our unified motion intermediate representation. Specifically, we employ D3DFR [5], which utilizes ResNet50 [13] to predict 3DMM coefficients. These coefficients consist of identity  $\alpha \in \mathbb{R}^{80}$ , expression  $\beta \in \mathbb{R}^{64}$ , texture  $\delta \in \mathbb{R}^{80}$ , illumination  $\gamma \in \mathbb{R}^{27}$ , and pose  $\rho \in \mathbb{R}^6$ . Therefore, given an input face  $I$ , we obtain the coefficient-based face descriptor  $P \in \mathbb{R}^{257}$ :

$$P = \psi^I(I) = \{\alpha, \beta, \delta, \gamma, \rho\}.$$

Given  $P$ , we acquire the reconstructed 3D face. By projecting it onto the 2D image plane using a fixed renderer  $R$ , we obtain the image-based face descriptor  $I_R$ :

$$I_R = \mathcal{R}(P).$$

### 6.2. Temporal Consistency

Following previous methods [14, 15, 50], we primarily focus on comparing the image quality and conducting an ablation study on the MIA module. To highlight the significance of the ICA module in generating smooth videos, we offer supplementary videos that qualitatively demonstrate the temporal stability of our proposed method alongside a quantitative evaluation presented in Tab. 2. Fréchet-Video Distance (FVD) [31, 35] is used to measure the temporal consistency of the generated videos. GAN-based methods are capable of producing smooth video results without explicit temporal modeling; however, due to their limited generative capacity, they may encounter challenges such as background blurring and facial distortions under large pose variations. In contrast, diffusion-based methods demonstrate excellent generalization and robustness, but their temporal stability is comparatively suboptimal. To maintain temporal stability, these methods often compromise a certain degree of motion control precision.

Methods	Architecture	FVD
MCNet [ICCV'23] [14]	GAN	616
HyperReenact [ICCV'23][2]	GAN	491
FADM [CVPR'23] [50]	Diffusion	285
AniPortrait [Arxiv'24] [38]	Diffusion	341
Ours w/o ICA	Diffusion	312
Ours	Diffusion	264

Table 2. Quantitative Evaluation on Temporal Consistency

In our approach, the MIA module delivers robust appearance features, ensuring temporal stability within clips,

while the ICA module effectively guarantees smooth transitions between clips. Most importantly, our method also ensures precise motion control, *i.e.*, pose, expression, and gaze.

### 6.3. More Results

**Handling Extreme Poses** In Fig. 10, we evaluate our approach on MPIE [10] to verify the capability to handle extreme poses. DG [16] is a method specifically designed to address large pose reenactment. Even without fine-tuning on MPIE, our method achieves better results than DG. This is because DG is constrained by the distribution of MPIE, and tends to alter the facial position within images and produces blurry results.

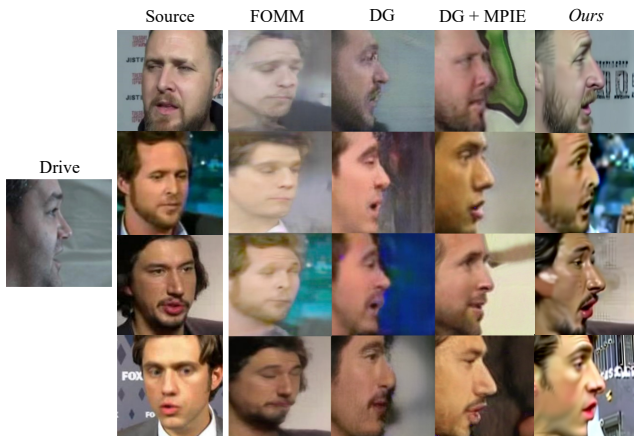


Figure 10. **Evaluation on MPIE.** We compare our method with DG using images directly cropped from the paper. "DG" represents results obtained by training DG on VoxCeleb, while "DG+MPIE" denotes results after fine-tuning DG on MPIE.

**Audio-driven Face Animation** Our model can be directly integrated with the audio2exp and audio2pose modules of SadTalker. In Fig. 12, given that our approach excels in generating fine-grained appearance details under large poses and offers added control over gaze, it yields results that are more vivid and clearer than those produced by the face render in SadTalker and AniPortrait. This demonstrates the advantage of our motion-modulated appearance features in generating subtle facial expressions.

**Intuitive Facial Editing** 3DMM-based methods inherently support free manipulation of pose and expression. Notable techniques include PIRenderer and StyleHEAT. Our method also enables the free control of gaze. Fig. 11-left displays our results in controlling pose, expression, and gaze. Our method decouples and controls various motions to a greater extent and can generate fine-grained appearance details, such as teeth and hair. Fig. 11-right Compared to previous methods, our approach maintains structure and produces clearer images under similar levels of pose control,

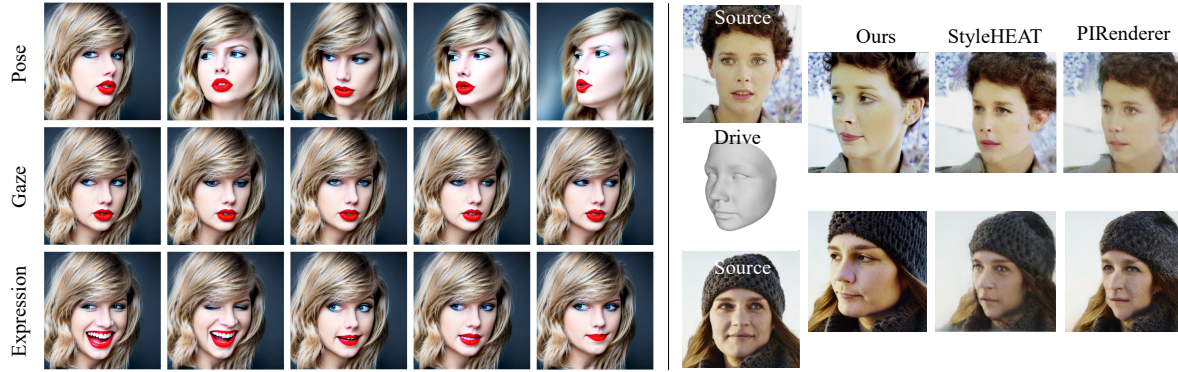


Figure 11. **Visualization of intuitive facial editing.** **Left:** Results in manipulating pose, expression, and gaze are shown. While 3DMM-based methods like PIRenderer and StyleHEAT inherently support pose and expression manipulation, our approach further offers free control of gaze. It exhibits an enhanced capability to decouple and control various motions and is adept at generating intricate appearance details like teeth and hair. **Right:** In comparison to prior methods, our method demonstrates superior structural preservation and clarity, especially under the same pose control, highlighting a marked elevation in the fidelity of facial expression nuances.



Figure 12. **Audio-driven Face Animation.** Our model is integrated with the audio2exp and audio2pose modules of SadTalker.

with a notable improvement in the facial expression details.

**Intriguing Real-World Applications** In Fig. 13, we showcase intriguing real-world applications. Thanks to the robust generalizability, high generative quality, and precise motion control of our method, we can leverage various amusing meme images from real life to animate specific real-world faces, offering strong entertainment value.

**Higher Resolution Results (768x768)** In the main paper, we utilize the VoxCeleb dataset with a resolution of 256 and 512, employing Stable Diffusion 1.5. Here, we showcase visual results obtained from the high-definition VFHQ [41] dataset at a resolution of 768. In Fig. 14, even when dealing with a profile source face, our method maintains a rational facial structure, in contrast to prior methods that often yield blurry results. This serves as validation for the efficacy of our approach in handling large-motion reenactment in high-

resolution scenarios.

**Cross-Domain Reenactment** In real-world scenarios, faces in the wild extend beyond authentic photographs of human faces to encompass a diverse array of artistic facial representations. To assess the generalizability of our approach in reenacting source faces across distinct domains, we conduct tests on faces featuring various artistic styles in Fig. 15. Our approach adeptly reenacts an array of faces, yielding high-fidelity outcomes. This success can be attributed to the pre-trained diffusion model, a capability beyond the reach of previous methods.

## References

- [1] Madhav Agarwal, Rudrabha Mukhopadhyay, Vinay P Nambodiri, and CV Jawahar. Audio-visual face reenactment. In





Figure 13. Intriguing Real-World Applications.

- Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5178–5187, 2023. 2
- [2] Stella Bounareli, Christos Tzelepis, Vasileios Argyriou, Ioannis Patras, and Georgios Tzimiropoulos. Hyperreenact: one-shot reenactment via jointly learning to refine and re-target faces. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7149–7159, 2023. 1, 9
- [3] Stella Bounareli, Christos Tzelepis, Vasileios Argyriou, Ioannis Patras, and Georgios Tzimiropoulos. One-shot neural face reenactment via finding directions in gan’s latent space. *International Journal of Computer Vision*, pages 1–31, 2024. 2
- [4] Di Chang, Yichun Shi, Quankai Gao, Jessica Fu, Hongyi Xu, Guoxian Song, Qing Yan, Xiao Yang, and Mohammad Soleymani. Magicdance: Realistic human dance video generation with motions & facial expressions transfer. *arXiv preprint arXiv:2311.12052*, 2023. 1, 3
- [5] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 0–0, 2019. 9
- [6] Chenpeng Du, Qi Chen, Tianyu He, Xu Tan, Xie Chen, Kai Yu, Sheng Zhao, and Jiang Bian. Dae-talker: High fidelity speech-driven talking face generation with diffusion auto-encoder. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 4281–4289, 2023. 3
- [7] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech-driven 3d facial animation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18770–18780, 2022. 3
- [8] Yue Gao, Yuan Zhou, Jinglu Wang, Xiao Li, Xiang Ming, and Yan Lu. High-fidelity and freely controllable talking head video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5609–5619, 2023. 2
- [9] Manuel B Garcia and Ahmed Mohamed Fahmy Yousef. Cognitive and affective effects of teachers’ annotations and talking heads on asynchronous video lectures in a web development course. *Research and Practice in Technology Enhanced Learning*, 18:020–020, 2023. 3
- [10] Ralph Gross, Iain Matthews, Jeffrey Cohn, Takeo Kanade, and Simon Baker. Multi-pie. *Image and vision computing*, 28(5):807–813, 2010. 9
- [11] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 1, 3
- [12] Yue Han, Jiangning Zhang, Junwei Zhu, Xiangtai Li, Yanhao Ge, Wei Li, Chengjie Wang, Yong Liu, Xiaoming Liu, and Ying Tai. A generalist face via learning unified facial representation. *arXiv preprint arXiv:2401.00551*, 2023. 2
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 9
- [14] Fa-Ting Hong and Dan Xu. Implicit identity representation conditioned memory compensation network for talking head video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23062–23072, 2023. 1, 6, 9
- [15] Fa-Ting Hong, Longhao Zhang, Li Shen, and Dan Xu. Depth-aware generative adversarial network for talking head video generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3397–3406, 2022. 1, 2, 9
- [16] Gee-Sern Hsu, Chun-Hung Tsai, and Hung-Yi Wu. Dual-generator face reenactment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 642–650, 2022. 9
- [17] Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. *arXiv preprint arXiv:2311.17117*, 2023. 1, 2, 3
- [18] Ricong Huang, Peiwen Lai, Yipeng Qin, and Guanbin Li. Parametric implicit face representation for audio-driven facial reenactment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12759–12768, 2023. 3
- [19] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. Curricularface: adaptive curriculum learning loss for deep face recognition. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5901–5910, 2020. 6
- [20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 6
- [21] Arsha Nagrani, Joon Son Chung, and Andrew Senior. Voxceleb: a large-scale speaker identification dataset. *arXiv preprint arXiv:1706.08612*, 2017. 5
- [22] Yuval Nirkin, Yosi Keller, and Tal Hassner. Fsgan: Subject agnostic face swapping and reenactment. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7184–7193, 2019. 1, 2
- [23] Gan Pei, Jiangning Zhang, Menghan Hu, Guangtao Zhai, Chengjie Wang, Zhenyu Zhang, Jian Yang, Chunhua Shen, and Dacheng Tao. Deepfake generation and detection: A benchmark and survey. *arXiv preprint arXiv:2403.17881*, 2024. 1
- [24] Xu Peng, Junwei Zhu, Boyuan Jiang, Ying Tai, Donghao Luo, Jiangning Zhang, Wei Lin, Taisong Jin, Chengjie Wang, and Rongrong Ji. Portraitbooth: A versatile portrait model for fast identity-preserved personalization. *arXiv preprint arXiv:2312.06354*, 2023. 2
- [25] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM international conference on multimedia*, pages 484–492, 2020. 3



- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. **1**
- [27] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13759–13768, 2021. **1, 2, 6, 7**
- [28] Shuai Shen, Wanhua Li, Zheng Zhu, Yueqi Duan, Jie Zhou, and Jiwen Lu. Learning dynamic facial radiance fields for few-shot talking head synthesis. In *European Conference on Computer Vision*, pages 666–682. Springer, 2022. **3**
- [29] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *Advances in neural information processing systems*, 32, 2019. **1, 2, 6, 7**
- [30] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2377–2386, 2019. **1, 2**
- [31] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3626–3636, 2022. **9**
- [32] Shuai Tan, Bin Ji, and Ye Pan. Flowvqtalker: High-quality emotional talking face generation through normalizing flow and quantization. *arXiv preprint arXiv:2403.06375*, 2024. **3**
- [33] Jiale Tao, Biao Wang, Borun Xu, Tiezheng Ge, Yuning Jiang, Wen Li, and Lixin Duan. Structure-aware motion transfer with deformable anchor model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3637–3646, 2022. **1, 6, 7**
- [34] Linrui Tian, Qi Wang, Bang Zhang, and Liefeng Bo. Emo: Emote portrait alive-generating expressive portrait videos with audio2video diffusion model under weak conditions. *arXiv preprint arXiv:2402.17485*, 2024. **1, 3**
- [35] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. **9**
- [36] Suzhen Wang, Yifeng Ma, Yu Ding, Zhipeng Hu, Changjie Fan, Tangjie Lv, Zhidong Deng, and Xin Yu. Styletalk++: A unified framework for controlling the speaking styles of talking heads. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. **3**
- [37] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10039–10049, 2021. **6, 7**
- [38] Huawei Wei, Zejun Yang, and Zhisheng Wang. Anipportrait: Audio-driven synthesis of photorealistic portrait animation. *arXiv preprint arXiv:2403.17694*, 2024. **1, 3, 5, 6, 9**
- [39] Jianzong Wu, Xiangtai Li, Chenyang Si, Shangchen Zhou, Jingkang Yang, Jiangning Zhang, Yining Li, Kai Chen, Yunhai Tong, Ziwei Liu, et al. Towards language-driven video inpainting via multimodal large language models. *arXiv preprint arXiv:2401.10226*, 2024.
- [40] Jianzong Wu, Xiangtai Li, Yanhong Zeng, Jiangning Zhang, Qianyu Zhou, Yining Li, Yunhai Tong, and Kai Chen. Motionbooth: Motion-aware customized text-to-video generation. *arXiv pre-print arXiv:2406.17758*, 2024. **3**
- [41] Liangbin Xie, Xintao Wang, Honglun Zhang, Chao Dong, and Ying Shan. Vfhq: A high-quality dataset and benchmark for video face super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 657–666, 2022. **10**
- [42] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven 3d facial animation with discrete motion prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12780–12790, 2023. **3**
- [43] Chao Xu, Jiangning Zhang, Yue Han, Guanzhong Tian, Xianfang Zeng, Ying Tai, Yabiao Wang, Chengjie Wang, and Yong Liu. Designing one unified framework for high-fidelity face reenactment and swapping. In *European Conference on Computer Vision*, pages 54–71. Springer, 2022. **1, 2**
- [44] Chao Xu, Jiangning Zhang, Miao Hua, Qian He, Zili Yi, and Yong Liu. Region-aware face swapping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7632–7641, 2022. **2**
- [45] Chao Xu, Shaoting Zhu, Junwei Zhu, Tianxin Huang, Jiangning Zhang, Ying Tai, and Yong Liu. Multimodal-driven talking face generation via a unified diffusion-based generator. *CoRR (2023)*, pages 1–14, 2023. **3**
- [46] Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation using diffusion model. *arXiv preprint arXiv:2311.16498*, 2023. **1, 2, 3, 5**
- [47] Kewei Yang, Kang Chen, Daoliang Guo, Song-Hai Zhang, Yuan-Chen Guo, and Weidong Zhang. Face2face  $p$ : Real-time high-resolution one-shot face reenactment. In *European conference on computer vision*, pages 55–71. Springer, 2022. **1, 2**
- [48] Zhenhui Ye, Tianyun Zhong, Yi Ren, Jiaqi Yang, Weichuang Li, Jiawei Huang, Ziyue Jiang, Jinzheng He, Rongjie Huang, Jinglin Liu, Chen Zhang, Xiang Yin, Zejun MA, and Zhou Zhao. Real3d-portrait: One-shot realistic 3d talking portrait synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. **3**
- [49] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan. In *European conference on computer vision*, pages 85–101. Springer, 2022. **1, 2, 3**
- [50] Bohan Zeng, Xuhui Liu, Sicheng Gao, Boyu Liu, Hong Li, Jianzhuang Liu, and Baochang Zhang. Face animation with an attribute-guided diffusion model. In *Proceedings of*

- the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 628–637, 2023. 2, 6, 7, 9
- [51] Bowen Zhang, Chenyang Qi, Pan Zhang, Bo Zhang, Hsiang-Tao Wu, Dong Chen, Qifeng Chen, Yong Wang, and Fang Wen. Metaportrait: Identity-preserving talking head generation with fast personalized adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22096–22105, 2023. 1, 2
- [52] Jiangning Zhang, Liang Liu, Zhucun Xue, and Yong Liu. Apb2face: Audio-guided face reenactment with auxiliary pose and blink signals. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4402–4406. IEEE, 2020. 3
- [53] Jiangning Zhang, Xianfang Zeng, Mengmeng Wang, Yusu Pan, Liang Liu, Yong Liu, Yu Ding, and Changjie Fan. Freenet: Multi-identity face reenactment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5326–5335, 2020. 1, 2
- [54] Jiangning Zhang, Xianfang Zeng, Chao Xu, and Yong Liu. Real-time audio-guided multi-face reenactment. *IEEE Signal Processing Letters*, 29:1–5, 2021. 3
- [55] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6
- [56] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8652–8661, 2023. 3
- [57] Yiming Zhang, Zhening Xing, Yanhong Zeng, Youqing Fang, and Kai Chen. Pia: Your personalized image animator via plug-and-play modules in text-to-image models. *arXiv preprint arXiv:2312.13964*, 2023. 3
- [58] Jian Zhao and Hui Zhang. Thin-plate spline motion model for image animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3657–3666, 2022. 2, 6, 7
- [59] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makelttalk: speaker-aware talking-head animation. *ACM Transactions On Graphics (TOG)*, 39(6):1–15, 2020. 3

Source  
Drive

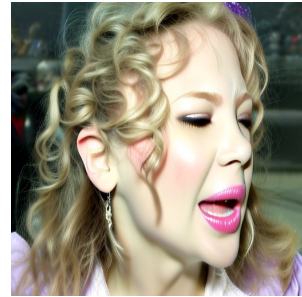
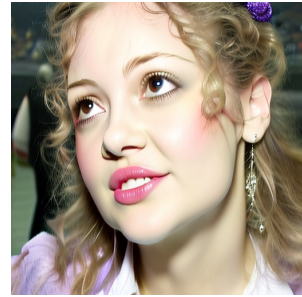
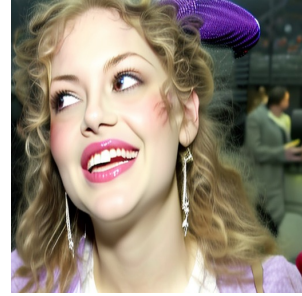
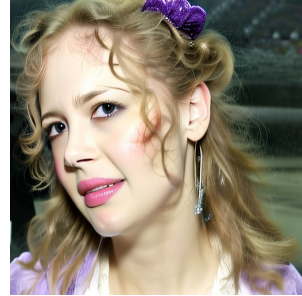
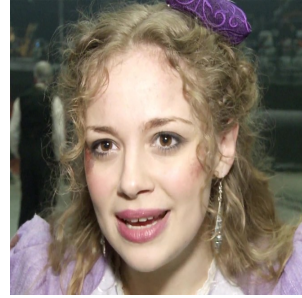


Figure 14. Higher Resolution Results (768\*768).

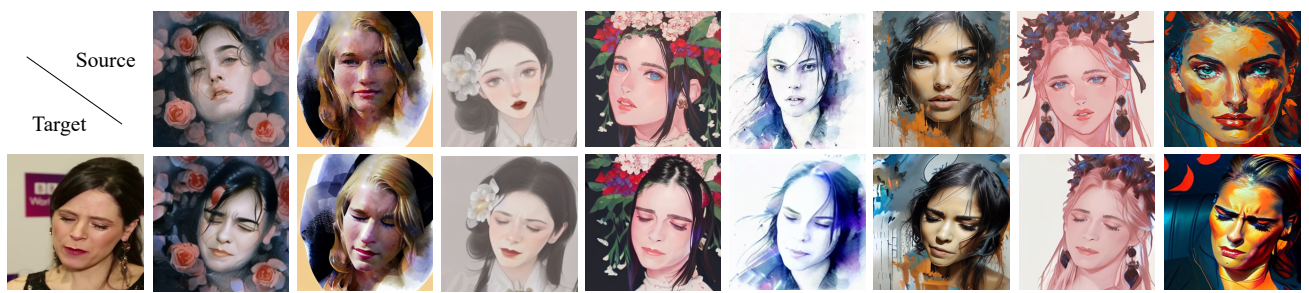


Figure 15. Cross-Domain Reenactment. Real faces effectively reenact artistic facial arts.