





Face-Adapter for Pre-trained Diffusion Models with Fine-Grained ID and Attribute Control

Yue Han¹, Junwei Zhu², Keke He², Xu Chen², Yanhao Ge³, Wei Li³,
Xiangtai Li⁴, Jiangning Zhang^{1,2}()^B, Chengjie Wang², and Yong Liu¹()^B

¹ Zhejiang University, Hangzhou, China
12432015@zju.edu.cn, yongliu@iipc.zju.edu.cn

² Tencent, Shenzhen, China

³ VIVO, Shanghai, China

⁴ Nanyang Technological University, Singapore, China
<https://faceadapter.github.io/face-adapter.github.io/>

Abstract. Current face reenactment and swapping methods mainly rely on GAN frameworks, but recent focus has shifted to pre-trained diffusion models for their superior generation capabilities. However, training these models is resource-intensive, and the results have not yet achieved satisfactory performance levels. To address this issue, we introduce **Face-Adapter**, an efficient and effective adapter designed for high-precision and high-fidelity face editing for pre-trained diffusion models. We observe that both face reenactment/swapping tasks essentially involve combinations of target structure, ID and attribute. We aim to sufficiently decouple the control of these factors to achieve both tasks in one model. Specifically, our method contains: 1) A Spatial Condition Generator that provides precise landmarks and background; 2) A Plug-and-play Identity Encoder that transfers face embeddings to the text space by a transformer decoder. 3) An Attribute Controller that integrates spatial conditions and detailed attributes. Face-Adapter achieves comparable or even superior performance in terms of motion control precision, ID retention capability, and generation quality compared to fully fine-tuned face reenactment/swapping models. Additionally, Face-Adapter seamlessly integrates with various StableDiffusion models.

Keywords: Face Reenactment · Face Swapping · Diffusion Model

1 Introduction

Face reenactment aims to transfer the target motion onto the source identity and attributes, while face swapping aims to transfer the source identity onto the target motion and attributes. Both tasks require complete disentangling and fine-grained control of identity, attributes, and motion. Current face reenactment and swapping techniques mainly rely on GAN-based frameworks [2, 3, 11, 15, 23, 26, 28, 29, 32, 47]. However, GAN-based methods encounter

Y. Han and J. Zhu—Co-first authors.

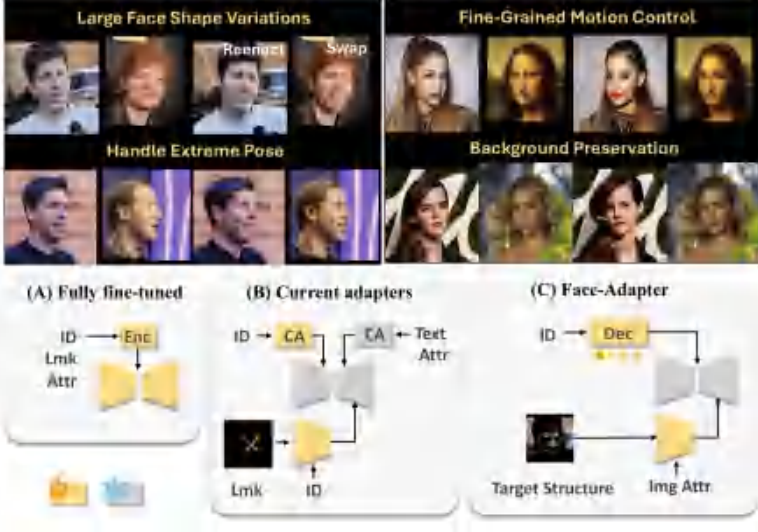


Fig. 1. Top: Face-Adapter supports a ‘one-model-two-tasks’ approach and demonstrates robustness under various challenging scenarios. **Bottom:** The design motivation is (1) Both face reenactment and swapping require fully disentangled ID, target structure, and attribute control; (2) Addressing overlooked issues unified in target structure; (3) Effective ID injection avoids SD fine-tuning, making Face-Adapter plug-and-play.

limitations in their generative capabilities, making it challenging to tackle hard cases, such as handling large poses in face reenactment and accommodating facial shape variations in face swapping.

Existing studies [41, 48] have attempted to address these challenges by leveraging the powerful generative capabilities of the diffusion models. However, these methods necessitate full model training, resulting in significant computational overhead, and they have not been successful in delivering satisfactory outcomes. For instance, FADM [41] refines the results of GAN-based reenactment methods, which improves image quality but still fails to resolve the blurring issue caused by large pose variation. On the other hand, DiffSwap [48] produces blurry facial outcomes due to the lack of background information during training, which hampers model learning. Moreover, these methods do not fully exploit the potential of large pre-trained diffusion models. To reduce training costs, some methods [30, 38] have introduced face editing adapter plugins for large pre-trained diffusion models. However, these approaches primarily focus on attribute editing using text, which inevitably weakens spatial control to ensure text editability. For example, they can only use five points [30] to control facial poses, limiting their ability to control expressions and gaze precisely. On the other hand, direct inpainting with masks of the face area does not take into account facial shape changes, leading to a decrease in identity preservation.

To address the above challenges, we are committed to developing an efficient and effective face editing adapter (**Face-Adapter**) for pre-trained diffusion

models, specifically targeting face reenactment and swapping tasks. The design motivation of Face-Adapter is threefold: (1) Fully disentangled ID, target structure, and attribute control enable a ‘one-model-two-tasks’ approach; (2) Addressing overlooked issues; (3) Simple yet effective, plug and play. Specifically, the proposed Face-Adapter comprises three components: **1)** Spatial Condition Generator (SCG in Sect. 3.1) is designed to automatically predict 3D prior landmarks and the mask of the varying foreground area, which provides more reasonable and precise guidance for subsequent controlled generation. In addition, for face reenactment, this strategy mitigates potential problems that could occur when only extracting the background from the source image, such as inconsistencies caused by alterations in the target background due to the movement of the camera or face objects; For face swapping, the model learns to maintain background consistency, glean clues about global lighting and spatial reference, and try to generate content in harmony with the background. **2)** Identity Encoder (IE in Sect. 3.2) uses the pre-trained recognition model to extract face embeddings and then transfers them to the text space by learnable queries from the transformer decoder. This manner greatly improves the identity consistency of the generated images. **3)** Attribute Controller (AC in Sect. 3.3) includes two sub-modules: The spatial control combines the landmarks of target motion with the unchanged background obtained from the Spatial Condition Generator. The attribute template supplements the absent attribute, encompassing lighting, a portion of the background, and hair. Both two tasks can be perceived as a procedure that executes conditional inpainting, utilizing the provided identity and absent attribute content. This process adheres to the stipulations of the given spatial control, attaining congruity and harmony with the background. Our contributions can be summarized as follows:

- We introduce Face-Adapter, a lightweight facial editing adapter designed to facilitate precise control over identity and attributes for pre-trained diffusion models. This adapter efficiently and proficiently tackles face reenactment and swapping tasks, surpassing previous state-of-the-art GAN-based and diffusion-based methods.
- We propose a novel Spatial Condition Generator module to predict the requisite generation areas, collaborating with the Identity Encoder and Attribute Controller to frame reenactment and swapping tasks as conditional inpainting with sufficient spatial guidance, identity, and essential attributes. Through reasonable and highly decoupled condition designs, we unleash the generative capabilities of pre-trained diffusion models for both tasks.
- Face-Adapter serves as a training-efficient, plug-and-play, face-specific adapter for pre-trained diffusion models. By freezing all parameters in the denoising U-Net, our method effectively capitalizes on priors and prevents overfitting. Furthermore, Face-Adapter supports a “one model for two tasks” approach, enabling simple input modifications to independently accomplish superior or competitive results of two facial tasks on VoxCeleb1/2 datasets.

2 Related Work

Face Reenactment involves extracting motion from a human face and transferring it to another face [1, 2, 21, 27, 34, 37, 42–44], which can be broadly divided into warping-based and 3DMM-based methods. *Warping-based methods* [13, 14, 27, 28, 31, 47] typically extract landmarks or region pairs to estimate motion fields and perform warping on the feature maps to transfer motions. When dealing with large motion variations, these methods tend to produce blurry and distorted results due to the difficulty in predicting accurate motion fields. *3DMM-based methods* [23] use facial reconstruction coefficients or rendered images from 3DMM as motion control conditions. The facial prior provided by 3DMM enables these methods to obtain more robust generation results in large pose scenarios. Despite offering accurate structure references, it only provides coarse facial texture and lacks references for hair, teeth, and eye movement. StyleHEAT [39] and HyperReenact [2] use StyleGAN2 to improve generation quality. However, StyleHEAT is limited by the dataset of frontal portraits, while HyperReenact suffers from resolution constraints and background blurring. To further improve generation quality, diffusion models have gained popularity. FADM [41] combines the previous reenactment model with diffusion refinements but the base model limits the driving accuracy. Recently, AnimateAnyone [16] employs heavy texture representation encoders (CLIP and a copy of U-Net) to ensure the textural quality of animated results, but this manner is costly. In contrast, we aim to leverage the generative capabilities of pre-trained T2I diffusion models fully and seek to comprehensively overcome the challenges presented in previous methods, *e.g.*, low -resolution generation, difficulty in handling large variations, efficient training, and unexpected artifacts.

Face Swapping aims to transfer the facial identity of the source image to the target image, with other attributes (*i.e.*, lighting, hair, background, and motion) of the target image unchanged. Recent methods can be broadly classified into GAN-based and diffusion-based approaches. **1)** Most GAN-based methods [3, 18, 19, 35, 36, 49] are dedicated to resolving the disentanglement and fusion of the identity and other attributes. Efforts include introducing face parsing masks, various losses for attribute-preserving, and designing fusion modules. Despite promising improvement, these methods often produce noticeable artifacts when dealing with significant changes in face shape or occlusions. HifiFace [32] alleviates this issue by utilizing 3DMM to reconstruct a reference face which combines the source face shape with other attributes of the target. However, relying on GAN to ensure generation quality, HifiFace still fails to inpaint harmonious results when dealing with large blank areas caused by face shape variation. **2)** Diffusion-based methods utilize the generative capabilities of the diffusion model to enhance sample quality. However, the numerous denoising steps during inference significantly increase the training costs when using attribute-preserving loss. DiffSwap [48] proposes midpoint estimation to address this issue, but the resulting error and the lack of background information for inpainting reference lead to unnatural results. Moreover, these methods require

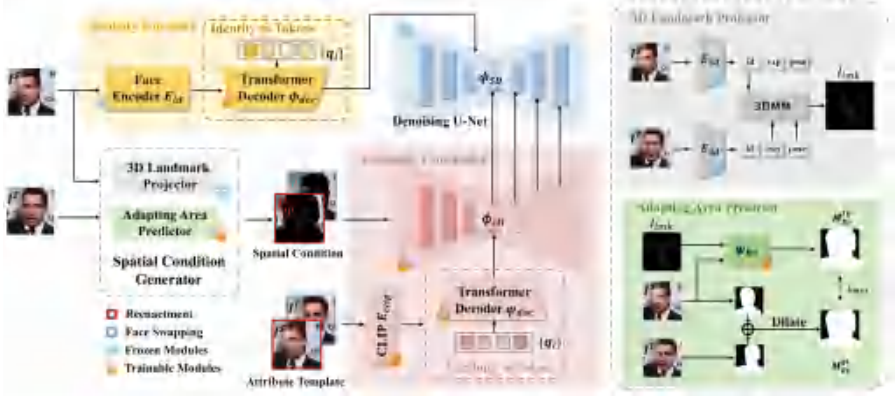


Fig. 2. Overview pipeline of our proposed Face-Adapter that consists of three modules: 1) The Spatial Condition Generator predicts 3D prior landmarks and adapts the foreground mask automatically, offering more accurate guidance for controlled generation. 2) The Identity Encoder improves identity consistency in generated images by transferring face embeddings to the text space using learnable queries. 3) The Attribute Controller features (i) spatial control that combines target motion landmarks with the invariant background from the Spatial Condition Generator, and (ii) an attribute template to fill in missing attributes.

costly training from scratch. In contrast, our Face-Adapter ensure image quality only relying on the denoise loss with complete disentanglement of the control of the target structure, ID and other attributes. Moreover, Face-Adapter further significantly reduces training costs by freezing all of U-Net’s parameters, which also preserves prior knowledge and prevents overfitting.

Personalization of Pretrained Diffusion Models. Personalization aims to insert a given identity into the pre-trained T2I diffusion models. Early works [10, 25] insert identity by using optimization or fine-tuning manners. Subsequent studies [4, 22, 33] introduce coarse spatial control, achieving multi-subject generation and regional attribute editing with text, but these methods require fine-tuning of most parameters. IP-adapter(-FaceID) [38] and InstantID [30] fine-tune only a few parameters. The latter achieves robust identity preservation. However, as a tradeoff for text editability, InstantID could only apply weak spatial control. Therefore, it struggles with fine movements (expression and gaze) in face reenactment and swapping. By comparison, our Face-Adapter is an effective and lightweight adapter designed for pre-trained diffusion models to accomplish face reenactment and swapping simultaneously.

3 Methods

The comprehensive structure of the proposed Face-Adapter is illustrated in Fig. 2, which aims to integrate identity into the attribute template, which provides essential attributes (*e.g.*, lighting, a portion of the background, and hair) based on the target motion (*e.g.*, pose, expression, and gaze).

3.1 Spatial Condition Generator

To provide more reasonable and precise guidance for subsequent controlled generation, we design a novel Spatial Condition Generator (SCG) to automatically predict 3D prior landmarks and the mask of the varying foreground area. In detail, this component consists of two sub-modules:

3D Landmark Projector. To surmount alterations in facial shape, we utilize a 3D facial reconstruction method [7] to extract the identity, expression individually and pose coefficients of the source and target faces. Subsequently, we recombine the identity coefficients of the source with the expression and pose coefficients of the target, reconstruct a new 3D face, and project it to acquire the corresponding landmarks.

Adapting Area Predictor. For face reenactment, prior methods assume that only the subject is in motion, while the background remains static in the training data. However, we observe that the background actually undergoes changes, encompassing the movement of both the camera and objects in the background, as illustrated in Fig. 3. If the model lacks knowledge of the background motion during training, it will learn to generate a blurry background. For face swapping, supplying the target background can also give the model clues about environmental lighting and spatial references. This added constraint of the background significantly diminishes the difficulty of the model learning, transitioning it from learning a task of generating from scratch to a task of conditional inpainting. As a result, the model becomes more attuned to preserving background consistency and generating content that seamlessly integrates with it.

In view of the above, we introduce a lightweight Adapting Area Predictor for both face reenactment and swapping, automatically predicting the region the model needs to generate (the adapting area) while maintaining the remaining area unchanged. For face reenactment, the adapting area constitutes the region occupied by the source image head before and after reenactment. We train a mask predictor φ_{Re} that accepts the target image I^T and motion landmarks I_{lmk} from the 3D Landmark Projector to predict the adapting area mask M_{Re}^{fg} . The mask ground truth M_{Re}^{gt} is generated by taking the union of the head regions (including hair, face, and neck) of the source and target, followed by outward



Fig. 3. Background inconsistency between the input (*i.e.*, source) and the groundtruth (*i.e.*, target) makes the model confused and fail to learn to generate clear background. Thus, we provide the background of the target image in the spatial condition during training to address this inconsistency.



Fig. 4. Comparisons with mask generated by pre-trained face parsing model (green) and φ_{Re} (white). The green mask cannot fully cover the entire portrait.

dilation. Head regions are obtained using a pre-trained face parsing model [40]. It should be noted that we cannot directly utilize the pre-trained face parsing model in face reenactment. As shown in Fig. 4 row 4, when the portrait area of the source image is larger (e.g., long hair and hat) than that in the target image, the green mask created by the pre-trained parsing model cannot fully cover the entire portrait and may result in artifacts at the boundary. However, the white mask created by φ_{Re} in Fig. 4 row 5 can encapsulate the whole portrait, as φ_{Re} merely uses the source image and 3D landmarks as input, and exhibits excellent generalization when the source and target images possess different identities.

For face-swapping, the adapting area constitutes the facial region of the target image I^T . We employ a pre-trained face parsing model [40] to predict the adapting area mask M_{Sw}^g of the target image I^T . Nonetheless, to accommodate face shape differences during testing, we designate the ground truth M_{Sw}^{gt} as the region obtained by dilating the facial area outward.

3.2 Identity Encoder

As demonstrated by IP-Adapter-FaceID [38] and InstantID [30], a high-level face embedding can ensure more robust identity preservation. As we observed, there is no need for heavy texture encoders [16] or additional identity networks [30] in face reenactment/swapping. By merely tuning a lightweight mapping module to map the face embedding into the fixed textual space, identity preservation is guaranteed. Specifically, given a face image I^S , the face embedding f_{id} is obtained by a pre-trained face recognition model E_{id} [6]. Subsequently, a three-layer transformer decoder ϕ_{dec} is employed to project the face embedding f_{id} into the fixed text semantic space of the pre-trained diffusion model, obtaining the identity tokens. The specified number N (we set $N = 77$ in this paper) of learnable queries $q_{id} = \{q_1, q_2, \dots, q_N\}$ in the transformer decoder constrains the sequence length of the identity embedding, ensuring it does not exceed the maximum length of the text embedding. Through this approach, the U-Net of the pre-trained diffusion model does not require any fine-tuning to adapt to the face embedding.

3.3 Attribute Controller

Spatial Control. In line with ControlNet [45], we create a copy of U-Net ϕ_{Ctl} and add spatial control I_{Sp} as the conditioning input. The spatial control image I_{Sp}^S/I_{Sp}^T is obtained by combining the target motion landmarks I_{Imk}^T and the non-adapting area obtained by the Adapting Area Predictor φ_{Re} (or φ_{Sw}).

$$I_{Sp}^S = I^S * (1 - M_{Re}^{fg}) + I_{Imk}^T, \text{ for face reenactment,}$$

$$I_{Sp}^T = I^T * (1 - M_{Sw}^{fg}) + I_{Imk}^T, \text{ for face swapping.}$$

At this juncture, both reenactment and swapping tasks can be viewed as processes of performing conditional inpainting, utilizing the given identity and other missing attribute content, following the provided spatial control.

Attribute Template. Given identity and spatial control with part of the background, the attribute template is designed to supplement the missing information, including lighting and part of the background and hair. Attribute embeddings $f_{attr} \in \mathbb{R}^{257*d}$ are extracted from the attribute template (I^S for reenactment and I^T for swapping) using CLIP E_{clip} . To simultaneously obtain local and global features, we use both the patch tokens and the global token. The feature mapper module is also constructed as a three-layer transformer layer φ_{dec} with learnable queries $q_{attr} = \{q_1, q_2, \dots, q_K\}$, $K = 77$.

3.4 Strategies for Boosting Performance

Training. 1) *Data Stream*: For both reenactment and face-swapping tasks, we use two images of the same person in different poses as source and target images. To support a “one model for both task” approach, we use a 50% probability to choose between reenactment and face-swapping data streams during training, i.e., the spatial control and attribute template in the Attribute Controller use the data streams indicated by red and blue respectively. 2) *Condition Dropping for Classifier-free Guidance*: The conditions we need to drop include identity tokens and attribute tokens input into the U-Net and ControlNet cross-attention. We use a 5% probability to simultaneously drop identity tokens and attribute conditions to enhance the realism of the image. To fully utilize the identity tokens for generation face images and improve identity preservation, we use an additional 45% probability to drop attribute tokens.

Inference. 1) *Adapting Area Predictor* : For reenactment, the input is the source (which is different from training) and corrected landmarks, and the output is the adapting area. For face-swapping, the input is the target, and the output is the adapting area. 2) *Negative Prompt for Classifier-Free Guidance*: For reenactment, negative prompts of both identity and attribute tokens are empty prompt embeddings. For face-swapping, to overcome the negative impact of the target identity in attribute tokens, we use the identity tokens of the target image as the negative prompt for identity tokens.

4 Experiments

4.1 Experimental Setup

Datasets. During training, we leverage the VoxCeleb1 and VoxCeleb2 [5] dataset. During the evaluation, we leverage the 491 test videos from the VoxCeleb1 [20] dataset and randomly sample 1,000 images in quantitative evaluation for face reenactment. We use FaceForensics++ [24] in quantitative evaluation for face swapping. We also spare 1,000 images from VoxCeleb2 for qualitative evaluation. Following the preprocessing method in FOMM [28], we crop faces from the original videos and resize them to 512×512 for training and evaluation.

Evaluation Metrics. For face reenactment, we use PSNR and LPIPS [46] to evaluate the reconstruction quality for same-identity reenactment. We use FID [12] to evaluate the overall quality of the generated images. We use cosine similarity (CSIM) calculated by [17] to evaluate identity preservation. The motion transfer error is measured by Pose, Exp, and Gaze, which calculate the average Euclidean distances of pose, expression, and gaze coefficients between the generated and drive images. For face swapping, ID retrieval (ID) retrieves the closest face to evaluate identity modification, while Pose, Exp, and Gaze evaluate the attribute error between the generated faces and target faces.

Implementation Details. The Adapting Area Predictor is modified from the parsing model [40], with 6 input channels and 1 output channel. The identity-to-tokens is implemented with a 3-layer transformer decoder, a linear layer is added to project the identity feature dimensions to 768. The architecture of attribute-to-tokens is the same as the identity-to-tokens, except the input dimensions of the linear layer are consistent with the output dimensions of the CLIP model. We adopt the StableDiffusion v1-5 [8] as the pre-trained diffusion model and clip-vit-large-patch14 [9] from OpenAI as the CLIP vision model in this paper. We train our face-adaptor for 70,000 steps on $8 \times V100$ NVIDIA GPUs with a constant learning rate of $1e-4$ and a batch size of 32.

4.2 Comparison with State-of-the-Art Methods

Face Reenactment. In Table 1, we compare with SoTA methods quantitatively on VoxCeleb1 test set, including GAN-based FOMM [28], PIRen-



Fig. 5. Same-identity face reenactment results on Voxceleb2 test set. Our method faithfully reconstructs the background and facial details.



Fig. 6. Cross-identity face reenactment results on Voxceleb2 test set.

derer [23], DG [15], TPSM [47], DAM [29], HyperReenact [2] and diffusion-based FADM [41]. FOMM, TPSM, and DAM are warping-based techniques, while PIRenderer and HyperReenact are 3DMM-based.

We achieve comparable or even optimal results in image quality. Owing to the Spatial Condition Generator, during training, incorporating the target background area in spatial condition avoids the interference of background motion. During inference, adding the source background in spatial condition significantly reduces the difficulty of generating backgrounds, improving background consistency. As a result, our method is capable of producing high-quality images with clear advantages in FID scores as well as in reconstruction metrics, *i.e.*, PSNR and LPIPS. In terms of motion control, our method performs well in pose and gaze error, but not as well in expression error. As our landmarks are derived from D3DFR, both the reconstruction and projection processes, along with the sparsity of the landmarks, result in a loss of expression accuracy. Therefore, our method achieves a relatively moderate performance in terms of expression error.

In Figs. 5 and 6, we compare with SoTA methods qualitatively on VoxCeleb1 and Voxceleb2 test set. The Spatial Condition Generator effectively ensures that our results are consistent with the source background and meanwhile reduces the training difficulty of the model, allowing it to focus more on face generation and improve the image quality. Freezing all parameters of the U-Net avoids overfitting and preserves as much of the powerful prior from the pre-trained diffusion model as possible. As a result, compared to other GAN-based methods and diffusion-based methods trained from scratch like FADM, our method is capable of generating faithful attribute details, *i.e.*, hair texture, hat, and accessories, that are consistent with the source image.



Fig. 7. Face swapping qualitative comparison results on Voxceleb2 test set.

In addition to local details, the attribute tokens in the Attribute Controller effectively extract global illumination from the source image, significantly outperforming other methods. This further highlights the strengths and capabilities of our proposed approach in capturing both local and global features, leading to more realistic and accurate results. Even when dealing with large poses, the Identity Encoder ensures robust identity preservation, and the pre-trained diffusion model reasonably generates attributes such as long hair that moves along with the face, demonstrating the superiority of our proposed adapter.

Face Swapping. In Table 2, we compare with SoTA methods quantitatively on FaceForensics++ test set, including GAN-based FaceShifter [18], SimSwap [3], HifiFace [32], InfoSwap [11], BlendFace [26] and diffusion-based DiffSwap [48].

Table 1. Quantitative evaluations among SoTAs on Voxceleb1 test set. **Bold** and underline correspond to the optimal and sub-optimal values, respectively.

Methods	Same-Identity					Cross-Identity						
	PSNR \uparrow	LPIPS \downarrow	FID \downarrow	Exp \downarrow	Pose \downarrow	Gaze \downarrow	CSIM \uparrow	Exp \downarrow	Pose \downarrow	Gaze \downarrow	CSIM \uparrow	FID \downarrow
FOMM [28]	22.77	<u>0.1344</u>	31.19	<u>2.92</u>	0.0276	<u>0.0566</u>	0.8499	6.89	0.0644	0.1003	0.539	51.57
PIRenderer [23]	21.65	0.1388	<u>29.98</u>	3.08	0.0409	0.0798	0.819	<u>6.42</u>	0.0646	0.0963	0.5361	40.71
DG [15]	14.01	0.4928	102.17	6.16	0.0707	0.112	0.0972	7.16	0.074	0.1287	0.0834	102.61
TPSM [47]	<u>23.8</u>	0.1367	34.11	2.70	0.0234	0.0627	0.8536	6.58	0.0548	0.0959	0.5514	54.83
DAM [29]	23.85	0.1484	38.6	2.87	0.027	0.0675	<u>0.8505</u>	6.82	0.0636	0.1034	0.5198	62.77
HyperReenact [2]	15.73	0.3361	88.72	3.68	0.0381	0.0743	0.5455	5.94	<u>0.0452</u>	<u>0.0812</u>	0.4665	88.02
FADM [41]	22.70	0.1392	31.58	3.11	0.0324	0.086	0.8472	7.03	0.0786	0.1239	<u>0.6152</u>	42.7
Ours	22.36	0.1281	29.27	3.24	<u>0.0243</u>	0.0415	0.7146	6.45	0.0355	0.0543	0.6429	<u>41.09</u>

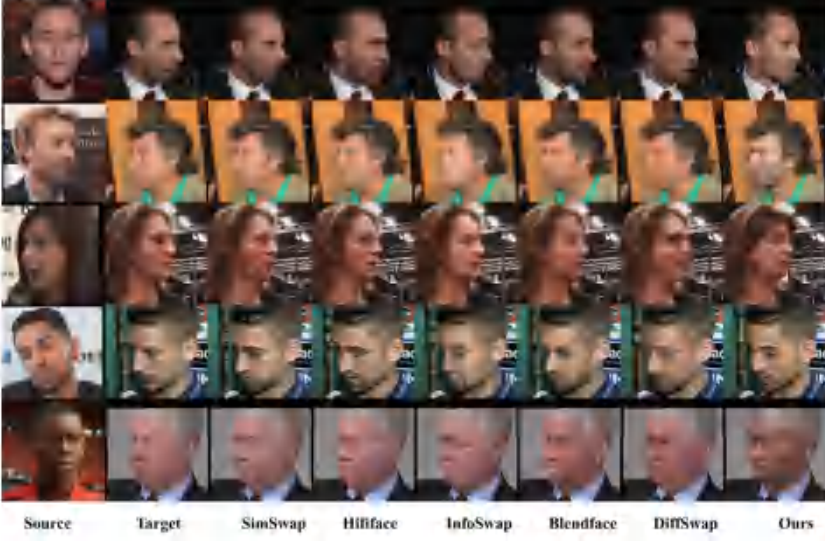


Fig. 8. Face swapping qualitative comparison results on Voxceleb2 test set.

Our 3D Landmark Projector helps to fuse the source face shape and target pose, expression and gaze to obtain the target motion landmarks in our spatial control. Our Adapting Area The predictor allows ample space for changes in face shape while keeping enough background for inpainting. This combined spatial condition benefits the model’s generation of natural images. Although DiffSwap also utilizes shape-aware landmarks via D3DFR as spatial control, its inpainting process only takes place during DDIM sampling. Lacking a background reference makes it difficult for the model to generate clear facial results, which significantly affects image quality and ID similarity. On the commonly used FaceForensics++ test set, our method is comparable to GAN-based methods in terms of ID, Pose, Exp and Gaze. Therefore, our method exhibits remarkable advantages in terms of ID while maintaining high motion accuracy compared to both GAN-based and diffusion-based SoTAs.

Figures 7 and 8 shows a qualitative comparison between our method and recent SoTA methods. Previous methods struggle with handling significant changes in face shape and large pose. When transferring a thin-faced person to a fat-faced target image, these methods typically maintain the face shape of the target image, leading to a significant loss of identity. In contrast, our spatial control effectively addresses the issue of face shape changes. Unlike previous approaches that merely crop out the facial region, our Adapting Area Predictor allows ample space for changes in face shape. With the powerful generation capability of the pre-trained SD model, we can naturally complete the regions with facial shape variations. Furthermore, by using the identity tokens of the target image as a negative prompt during face-swapping inference, we further enhance

Table 2. Quantitative results on the task of face swapping on FF++. Compared to the diffusion-based DiffSwap, our method significantly improves the metrics and achieves highly competitive results. **Note that our method can simultaneously perform both face reenactment and swapping.** **Bold** corresponds to the optimal values. *: evaluated results are from the official code. †: evaluated results are from the officially released generated videos.

Methods	ID ↑	Pose↓	Exp↓	Gaze↓
FaceShifter [18]†	87.99	0.0342	6.32	0.072
SimSwap [3]*	96.78	0.0261	5.94	0.0549
HifiFace [32]†	94.26	0.0382	6.50	0.0573
InfoSwap [11]*	99.26	0.0371	7.25	0.0617
BlendFace [26]*	89.91	0.0286	6.15	0.0556
DiffSwap [48]*	19.16	0.0237	4.94	0.0665
Ours	96.47	0.0319	6.66	0.0607

the identity similarity with the source face. As for large poses, previous methods struggle to generate plausible results, while our method directly generates faces from 3D landmarks without being affected by the pose.

4.3 Ablation Study and Further Analysis

We conducted an ablation study on the Adapting Area Predictor and assessed the necessity of fine-tuning CLIP. For a fair comparison, all three models here were trained for 35,000 steps. Quantitative evaluations are conducted on Voxceleb1 cross-identity test set for both face reenactment and swapping tasks.

Adapting Area Predictor. As demonstrated in Table 3 and Fig. 9, without the Adapting Area Predictor, the spatial control lacks the background and only includes landmarks from the 3D Landmark Projector. During training, the model extracts the background features from the source image in face reenactment, while using the target image background as the ground truth. This discrepancy tends to result in the model hallucinating background, and the model struggles to maintain consistency with the background of the source image during inference. As for face swapping, the model is not trained with inpainting task, which leads to noticeable unnatural artifacts when blending the face with the surrounding area during inference.

Table 3. Quantitative comparison of our model under different ablative configurations.

Methods	Face Reenactment					Face Swapping				
	FID↓	Pose↓	Exp↓	Gaze↓	ID↑	FID↓	Pose↓	Exp↓	Gaze↓	ID↑
w/o AAP	33.61	0.0281	3.72	0.045	0.6355	33.97	0.0395	6.13	0.0548	0.4530
w/o CLIP FT	33.09	0.0287	3.74	0.0435	0.6474	31.97	0.0396	6.21	0.0540	0.4696
Full Model	31.18	0.0266	3.61	0.0422	0.6616	30.78	0.0406	6.14	0.0547	0.4688



Fig. 9. Ablation study for Spatial Condition Generator and CLIP finetuning. The red boxes highlight the artifacts in the picture.

Fine-Tuning CLIP for Extracting Attribute Features. As demonstrated in Table 3 and Fig. 9, freezing the CLIP results in a decline in detailed attributes and image quality. The pre-trained CLIP is trained for discrimination tasks and lacks detailed texture features needed for generation tasks. Fine-tuning CLIP helps to extract detailed attribute features, including hair, clothing, part of the missing backgrounds, and global lighting; in addition to this, the fine-tuned CLIP model also extracts some features related to face identity, which benefits the identity similarity score in face reenactment.

5 Conclusion

In this paper, we present a novel Face-Adapter framework, a plug-and-play facial editing adapter that supports fine control over identity and attributes for pretrained diffusion models. Utilizing only one model, this adapter effectively addresses face reenactment and swapping tasks, surpassing previous state-of-the-art GAN-based and diffusion-based methods. Extensive qualitative and quantitative experiments demonstrate the superiority of our method.

Limitations. Our unified model is unable to achieve temporal stability in video face reenactment/ swapping, which requires incorporating additional temporal fine-tuning in the future.

Potential Social Impact. For the first time, we explore a lightweight framework based on diffusion for simultaneous face reenactment and swapping, which has higher practical value while improving the quality of generated content. However, the potential misuse of Face-Adapter can lead to privacy invasion, misinformation spread, and ethical concerns. To mitigate these risks, both visible and invisible digital watermarks can be incorporated to help identify the origin and authenticity of the content. On the other side, Face-Adapter can contribute to the field of forgery detection, further enhancing the ability to identify and combat deepfakes.

References

1. Agarwal, M., Mukhopadhyay, R., Namboodiri, V.P., Jawahar, C.: Audio-visual face reenactment. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 5178–5187 (2023)
2. Bounareli, S., Tzelepis, C., Argyriou, V., Patras, I., Tzimiropoulos, G.: Hyperreenact: one-shot reenactment via jointly learning to refine and retarget faces. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7149–7159 (2023)
3. Chen, R., Chen, X., Ni, B., Ge, Y.: Simswap: an efficient framework for high fidelity face swapping. In: Proceedings of the 28th ACM International Conference on Multimedia, pp. 2003–2011 (2020)
4. Choi, J., Choi, Y., Kim, Y., Kim, J., Yoon, S.: Custom-edit: Text-guided image editing with customized diffusion models. arXiv preprint [arXiv:2305.15779](https://arxiv.org/abs/2305.15779) (2023)
5. Chung, J.S., Nagrani, A., Zisserman, A.: Voxceleb2: Deep speaker recognition. arXiv preprint [arXiv:1806.05622](https://arxiv.org/abs/1806.05622) (2018)
6. Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: additive angular margin loss for deep face recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4690–4699 (2019)
7. Deng, Y., Yang, J., Xu, S., Chen, D., Jia, Y., Tong, X.: 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 (2019)
8. Face, H.: Runwayml stable diffusion v1.5. <https://huggingface.co/runwayml/stable-diffusion-v1-5>, Accessed on: yyyy-mm-dd
9. Foundations, M.: Openclip: Open-source implementation of clip (2022). https://github.com/mlfoundations/open_clip, Accessed on: yyyy-mm-dd
10. Gal, R., et al.: An image is worth one word: Personalizing text-to-image generation using textual inversion. arXiv preprint [arXiv:2208.01618](https://arxiv.org/abs/2208.01618) (2022)
11. Gao, G., Huang, H., Fu, C., Li, Z., He, R.: Information bottleneck disentanglement for identity swapping. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3404–3413 (2021)
12. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Adv. Neural Inform. Process. Syst. **30** (2017)
13. Hong, F.T., Xu, D.: Implicit identity representation conditioned memory compensation network for talking head video generation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 23062–23072 (2023)
14. Hong, F.T., Zhang, L., Shen, L., Xu, D.: Depth-aware generative adversarial network for talking head video generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3397–3406 (2022)
15. Hsu, G.S., Tsai, C.H., Wu, H.Y.: Dual-generator face reenactment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 642–650 (2022)
16. Hu, L., Gao, X., Zhang, P., Sun, K., Zhang, B., Bo, L.: Animate anyone: Consistent and controllable image-to-video synthesis for character animation. arXiv preprint [arXiv:2311.17117](https://arxiv.org/abs/2311.17117) (2023)
17. Huang, Y., et al.: Curricularface: adaptive curriculum learning loss for deep face recognition. In: proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5901–5910 (2020)

18. Li, L., Bao, J., Yang, H., Chen, D., Wen, F.: Faceshifter: Towards high fidelity and occlusion aware face swapping. arXiv preprint [arXiv:1912.13457](https://arxiv.org/abs/1912.13457) (2019)
19. Liu, Z., et al.: Fine-grained face swapping via regional gan inversion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8578–8587 (2023)
20. Nagrani, A., Chung, J.S., Zisserman, A.: Voxceleb: a large-scale speaker identification dataset. arXiv preprint [arXiv:1706.08612](https://arxiv.org/abs/1706.08612) (2017)
21. Nirkin, Y., Keller, Y., Hassner, T.: Fsgan: subject agnostic face swapping and reenactment. In: Proceedings of the IEEE/CVF international conference on computer vision, pp. 7184–7193 (2019)
22. Peng, X., et al.: Portraitbooth: A versatile portrait model for fast identity-preserved personalization. arXiv preprint [arXiv:2312.06354](https://arxiv.org/abs/2312.06354) (2023)
23. Ren, Y., Li, G., Chen, Y., Li, T.H., Liu, S.: Pirenderer: controllable portrait image generation via semantic neural rendering. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 13759–13768 (2021)
24. Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Nießner, M.: Faceforensics++: learning to detect manipulated facial images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1–11 (2019)
25. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: Dreambooth: fine tuning text-to-image diffusion models for subject-driven generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 22500–22510 (2023)
26. Shiohara, K., Yang, X., Taketomi, T.: Blendface: re-designing identity encoders for face-swapping. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7634–7644 (2023)
27. Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., Sebe, N.: Animating arbitrary objects via deep motion transfer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2377–2386 (2019)
28. Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., Sebe, N.: First order motion model for image animation. Adv. Neural Inform. Process. Syst. **32** (2019)
29. Tao, J., et al.: Structure-aware motion transfer with deformable anchor model. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3637–3646 (2022)
30. Wang, Q., Bai, X., Wang, H., Qin, Z., Chen, A.: Instantid: Zero-shot identity-preserving generation in seconds. arXiv preprint [arXiv:2401.07519](https://arxiv.org/abs/2401.07519) (2024)
31. Wang, T.C., Mallya, A., Liu, M.Y.: One-shot free-view neural talking-head synthesis for video conferencing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10039–10049 (2021)
32. Wang, Y., et al.: Hiface: 3d shape and semantic prior guided high fidelity face swapping. arXiv preprint [arXiv:2106.09965](https://arxiv.org/abs/2106.09965) (2021)
33. Xiao, G., Yin, T., Freeman, W.T., Durand, F., Han, S.: Fastcomposer: Tuning-free multi-subject image generation with localized attention. arXiv preprint [arXiv:2305.10431](https://arxiv.org/abs/2305.10431) (2023)
34. Xu, C., et al.: Designing one unified framework for high-fidelity face reenactment and swapping. In: European Conference on Computer Vision, pp. 54–71. Springer (2022). https://doi.org/10.1007/978-3-031-19784-0_4
35. Xu, C., Zhang, J., Hua, M., He, Q., Yi, Z., Liu, Y.: Region-aware face swapping. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7632–7641 (2022)

36. Xu, Z., Hong, Z., Ding, C., Zhu, Z., Han, J., Liu, J., Ding, E.: Mobilefaceswap: a lightweight framework for video face swapping. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 36, pp. 2973–2981 (2022)
37. Yang, K., Chen, K., Guo, D., Zhang, S.H., Guo, Y.C., Zhang, W.: Face2face ρ : Real-time high-resolution one-shot face reenactment. In: European Conference on Computer Vision, pp. 55–71. Springer (2022). https://doi.org/10.1007/978-3-031-19778-9_4
38. Ye, H., Zhang, J., Liu, S., Han, X., Yang, W.: Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint [arXiv:2308.06721](https://arxiv.org/abs/2308.06721) (2023)
39. Yin, F., et al.: Styleheat: one-shot high-resolution editable talking face generation via pre-trained stylegan. In: European Conference on Computer Vision, pp. 85–101. Springer (2022). https://doi.org/10.1007/978-3-031-19790-1_6
40. Yu, C., Wang, J., Peng, C., Gao, C., Yu, G., Sang, N.: BiSeNet: bilateral segmentation network for real-time semantic segmentation. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11217, pp. 334–349. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01261-8_20
41. Zeng, B., Liu, X., Gao, S., Liu, B., Li, H., Liu, J., Zhang, B.: Face animation with an attribute-guided diffusion model. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 628–637 (2023)
42. Zeng, X., Pan, Y., Wang, M., Zhang, J., Liu, Y.: Realistic face reenactment via self-supervised disentangling of identity and pose. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, pp. 12757–12764 (2020)
43. Zhang, B., et al.: Metaportrait: identity-preserving talking head generation with fast personalized adaptation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 22096–22105 (2023)
44. Zhang, J., et al.: Freenet: Multi-identity face reenactment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5326–5335 (2020)
45. Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 3836–3847 (2023)
46. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 586–595 (2018)
47. Zhao, J., Zhang, H.: Thin-plate spline motion model for image animation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3657–3666 (2022)
48. Zhao, W., Rao, Y., Shi, W., Liu, Z., Zhou, J., Lu, J.: Diffswap: high-fidelity and controllable face swapping via 3d-aware masked diffusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8568–8577 (2023)
49. Zhu, Y., Li, Q., Wang, J., Xu, C.Z., Sun, Z.: One shot face swapping on megapixels. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4834–4844 (2021)