

DiffFace: Diffusion-based Face Swapping with Facial Guidance

Kihong Kim^{*1} Yunho Kim^{*1} Seokju Cho² Junyoung Seo²
 Jisu Nam² Kychul Lee¹ Seungryong Kim^{†,2} Kwang Hee Lee^{†,1}

¹VIVE STUDIOS ²Korea University

¹{hxngiee, youknowyunho, lkc880425, lucas}@vivestudios.com

²{seokju_cho, se780, jisu_nam, seungryong_kim}@korea.ac.kr

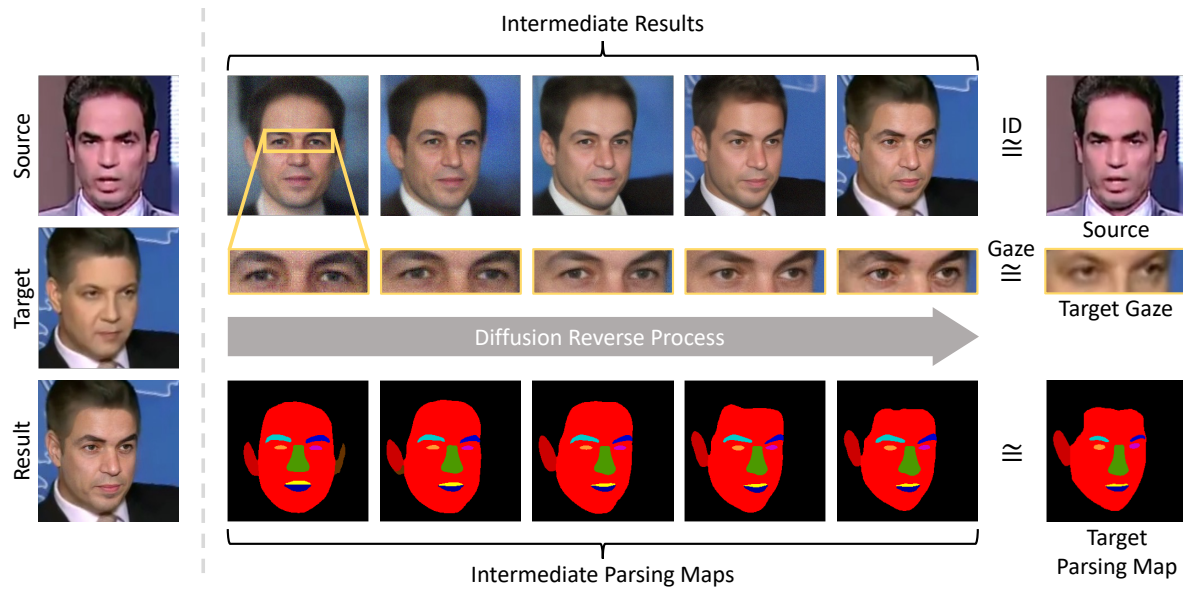


Figure 1. **Visualization of our novel diffusion-based face swapping framework, called DiffFace.** DiffFace gradually produces images with source identity and target attributes such as gaze, structure and pose.

Abstract

In this paper, we propose a novel diffusion-based face swapping framework, called DiffFace, composed of training ID Conditional DDPM, sampling with facial guidance, and a target-preserving blending.¹ In specific, in the training process, the ID Conditional DDPM is trained to generate face images with the desired identity. In the sampling process, we use the off-the-shelf facial expert models to make the model transfer source identity while preserving target attributes faithfully. During this process, to preserve the background of the target image, we additionally propose a target-preserving blending strategy. It helps our model to

keep the attributes of the target face from noise while transferring the source facial identity. In addition, without any re-training, our model can flexibly apply additional facial guidance and adaptively control the ID-attributes trade-off to achieve the desired results. To the best of our knowledge, this is the first approach that applies the diffusion model in face swapping task. Compared with previous GAN-based approaches, by taking advantage of the diffusion model for the face swapping task, DiffFace achieves better benefits such as training stability, high fidelity, and controllability. Extensive experiments show that our DiffFace is comparable or superior to the state-of-the-art methods on the standard face swapping benchmark.

^{*} Authors contributed equally.

[†] Corresponding author.

¹ Project Page : <https://hxngiee.github.io/DiffFace/>

1. Introduction

Generative adversarial networks (GANs) [13] have shown incredible empirical success in image generation [5, 16, 18, 21, 31, 59] tasks. Recently, most of face swapping tasks [1, 8, 12, 24, 52, 60] have been proposed based on GANs. Despite their empirical successes, it has been well known that the training of GANs is inherently unstable [32] due to the min-max optimization problem of the generator and discriminator. To alleviate this problem, complex architecture, various loss functions, and extensive hyperparameter tuning have been required. In addition, tuning hyperparameters becomes even more challenging if external models are combined.

Face swap is a task to synthesize an image with the identity of the source image while preserving the target image’s attributes (e.g., expression, pose, and shape). Off-the-shelf expert models trained on specific purposes can be used in the face swap task. For instance, ID embedder can be utilized to constrain the synthesized image’s identity to follow the source’s identity [8, 24]. Nevertheless, balancing identity and attributes is one of the most challenging problems. GAN-based face swapping tasks [8, 52] maintain this balance through a combination of loss functions related with ID and facial attributes. However, it is necessary to perform multiple training to find the desired hyperparameters. Due to the trade-off between ID and attributes, a greater focus on maintaining natural attributes often makes it difficult to achieve satisfactory results in transferring the source face ID to a synthesized face.

Recently, diffusion models [11, 34, 39, 42] have attracted much attention as an alternative to GANs [13]. Diffusion models, as opposed to GANs, enable more stable training, showing desirable results in terms of diversity and fidelity. To tackle the trade-off between fidelity and diversity, classifier guidance [11] is introduced to guide the diffusion model. Such a guidance technique is also widely used in conditional generation [27, 58], especially in text-to-image generation [3, 33, 41, 47]. Despite various advantages of the guidance technique, we stress the usability of the external module as guidance at test time.

In this paper, we propose a novel diffusion-based face swap framework, named *DiffFace*, which is composed of training ID Conditional DDPM, sampling with facial guidance, and a target-preserving blending strategy. First of all, we present an ID Conditional DDPM. In the training process, the ID Conditional DDPM is trained to generate face images with the desired identity. We make the diffusion model aware of facial identity by not only injecting the identity feature vector from the ID embedder but also posing additional constraint with the identity similarity loss. In the sampling process, we use facial guidance driven by various pretrained experts to enable the model to transfer source identity while preserving target attributes faithfully as illustrated in Fig. 1. During this process, to preserve the background of the tar-

get image and obtain the desired face swapping result, we additionally propose a target-preserving blending strategy. By gradually increasing the facial mask intensity over the time of the diffusion process, it prevents our model from completely forgetting the attributes of the target face by noise while transferring the source facial identity. In addition, our model can flexibly apply various facial guidance and adaptively control the ID-attributes trade-off to achieve the desired results without any re-training. Compared with the GAN-based face swapping works, *DiffFace* achieves better benefits such as training stability, high fidelity, and controllability. Experiments demonstrate that our method surpasses other state-of-the-art methods on standard face swapping benchmark.

2. Related Work

2.1. Diffusion Model

Diffusion models [15, 34] have attained much attention as a generative model showing desirable qualities while maintaining a higher distribution coverage. The sampling process of the diffusion models can be intuitively viewed as a denoising process. Utilizing the characteristics of the gradual denoising process, various works [29, 30, 46, 49, 50] have achieved remarkable results in the field of conditional generation [9], local image editing [29, 30, 50], and image translation [30, 46, 49, 50]. On the one hand, ADM [11] proposed to interpret the gradient of the external classifier as guidance and inject it during the reverse process. This not only improved the sample quality but also succeeded in conditional generation. These guidance techniques have been widely used, including text-to-image generation [33, 41, 42, 47]. Also, recent work [3] enabled text-guided editing of local area by injecting CLIP guidance [40] only into the local area. Although various methods have been proposed, no research directly tackled the face swapping task. For the first time, we propose a high-fidelity and controllable face swapping framework based on the conditional diffusion model.

2.2. Face Swap Model

Structural Prior-Guided Models. Previous face swapping methods such as HifiFace [52] combine structural information extracted from 3D Morphable Model [4] with GANs to produce 3D shape-aware identity. With improving the geometric structure of generated images, FSGAN [35] attempts to reenact the source image to match the target image using facial key points and segmentation. However, inaccurate structural or facial key points information makes previous models struggle to generate high-fidelity results.

Reconstruction-Based Models. DeepFakes [1] was trained to swap faces between paired identities. However, this method can only be applied to one specific identity. Subject-agnostic face swap methods such as Faceshifter [25]

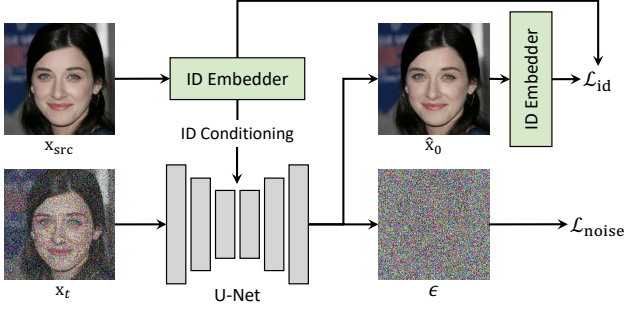


Figure 2. **Training procedure for ID Conditional DDPM.** $\hat{\mathbf{x}}_0$ denotes the predicted denoised image at timestep t . In each denoising step, we condition ID vector of the source image into the U-Net. Then, we compute cosine similarity loss between the ID vectors of source and $\hat{\mathbf{x}}_0$.

and SimSwap [8] were proposed to overcome the limitation of identity-specific face swap. These subject-agnostic models typically tune intermediate features of the target image to inject the identity of the source image. Even though these methods could synthesize arbitrary identities, they are insufficient to generate high fidelity results.

StyleGAN-Based Models. Recently, StyleGAN [21, 22] based face swap model [54] has emerged as a solution for high resolution face swap. These models relied on StyleGAN architectures [20–22], which have powerful latent space representation and expressibility. MegaFS [60] proposed to invert source and target image into latent space, then fuses these two features appropriately and feeds it into StyleGAN generator to obtain the results. Other work [53] focused more on separating identity and pose information. InfoSwap [12] proposed identity contrastive loss that better disentangles StyleGAN latent space. The StyleGAN-based face swap methods [12, 53, 54, 60] fusing the source and target’s multi-resolution features typically failed to synthesize fine details.

3. Preliminaries: Denoising Diffusion Probabilistic Models

Diffusion models generate a realistic image from a standard Gaussian distribution by reversing a recurrent noising process [15]. The forward process gradually alters to Gaussian distribution from the data $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, which is defined as $q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$, where β_t is a predefined variance schedule. In addition, the reverse process is as follows:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_\theta(\mathbf{x}_t, t)\mathbf{I}), \quad (1)$$

where $\mu_\theta(\cdot)$ and $\sigma_\theta(\cdot)$ can be parameterized using deep neural networks. But, in practice, it is known that using noise approximation model $\epsilon_\theta(\mathbf{x}_t, t)$ worked best instead of using

Algorithm 1 Training ID Conditional DDPM

```

repeat
   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
   $t \sim \text{Uniform}(\{1, \dots, T\})$ 
   $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
   $\mathbf{x}_t \sim \mathcal{N}(\sqrt{\alpha_t}\mathbf{x}_0, (1 - \alpha_t)\mathbf{I})$ 
   $\mathbf{v}_{\text{id}} \leftarrow \mathcal{D}_t(\mathbf{x}_0)$ 
   $\hat{\mathbf{x}}_0 \leftarrow \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t}\epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_{\text{id}})}{\sqrt{\alpha_t}}$ 
   $\hat{\mathbf{v}}_{\text{id}} \leftarrow \mathcal{D}_t(\hat{\mathbf{x}}_0)$ 
  Take gradient descent step on
   $\nabla_\theta (\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_{\text{id}})\|_2^2 + \lambda \|\mathbf{v}_{\text{id}} - \hat{\mathbf{v}}_{\text{id}}\|_2^2)$ 
until converged

```

$\mu_\theta(\mathbf{x}_t, t)$. [15] Thus, $\mu_\theta(\mathbf{x}_t, t)$ can be induced as:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(\mathbf{x}_t, t) \right). \quad (2)$$

Given \mathbf{x}_t , the reverse process of the diffusion model usually outputs \mathbf{x}_{t-1} . But, we can also directly derive $\hat{\mathbf{x}}_0$ which is the fully denoised prediction given \mathbf{x}_t by $\hat{\mathbf{x}}_0 = f_\theta(\mathbf{x}_t, t)$ where

$$f_\theta(\mathbf{x}_t, t) := \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t}\epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}}. \quad (3)$$

We utilize this fully denoised prediction $\hat{\mathbf{x}}_0$ for facial expert modules.

Meanwhile, ADM [41] proposes a guidance technique for the diffusion model. ADM [41] trains a classifier $p(y|\mathbf{x}_t, t)$, which takes a noised image as input, and regards the gradient of the classifier as a guidance for the diffusion models:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, y) := \mathcal{N}(\mu + s\nabla_{\mathbf{x}_t} p(y|\mathbf{x}_t, t), \sigma\mathbf{I}), \quad (4)$$

where s is a constant for the guidance scale and μ, σ are $\mu_\theta(\mathbf{x}_t, t)$ and $\sigma_\theta(\mathbf{x}_t, t)$, respectively. Note that the guidance technique proposed by ADM [41] uses an unconditional diffusion model, which is different in that our work utilizes ID Conditional DDPM. Although the diffusion model has the advantage of adding various guidance techniques while guaranteeing strong generation capabilities, employing the diffusion model to face swapping has not yet been explored due to the difficulties described below.

4. Methodology

In this section, we describe our method called *DiffFace*. The whole process is divided into training ID Conditional DDPM, sampling with facial guidance, and target-preserving blending. We propose ID Conditional DDPM which constructs diffusion model suitable for face swap tasks. Then, we describe the facial guidance designed to synthesize the desired image in the diffusion process. Lastly, we introduce target-preserving blending, which can help ID Conditional DDPM preserve the target facial detail better.

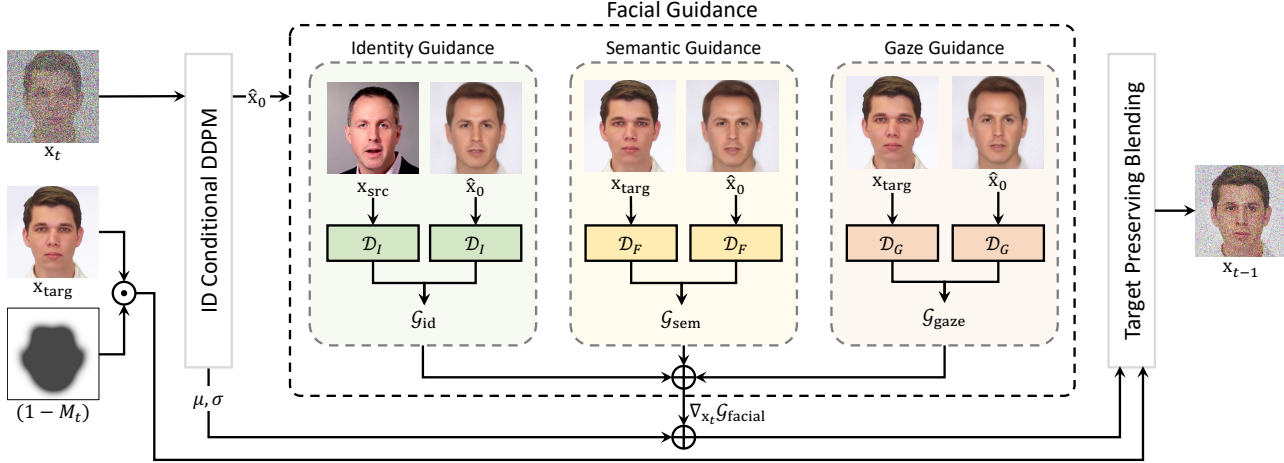


Figure 3. **Sampling procedure with facial guidance.** Given facial mask and target image, we leverage the generative process of diffusion model to synthesize coherent face swapping results. Facial guidance induces the masked region to have the desired facial attributes. Our facial guidance module is composed of three components, which can be found in the bottom of the figure.

Throughout this paper, we use \mathbf{x}_{src} , \mathbf{x}_{targ} , and \mathbf{x}_{swap} to represent the source, target, and synthesized image, respectively. Also, let us denote \mathcal{D} by the pretrained guidance network. Specifically, \mathcal{D}_I denotes identity embedder [10, 51], \mathcal{D}_F denotes face parser [55], and \mathcal{D}_G denotes gaze estimator [36].

4.1. ID Conditional DDPM

Our key idea for face swapping with the diffusion model is to inject the identity feature into the diffusion model. While previous methods have studied these conditioning problems [41, 42, 46, 47] extensively, no study infuses identity information as a condition into the diffusion model. Thus, we employ the structure of the conditional diffusion model, where additional information can be injected. As shown in Fig. 2 we first inject the source image \mathbf{x}_{src} in the identity embedder \mathcal{D}_I (e.g., ArcFace [10] and CosFace [51]) to obtain the source identity \mathbf{v}_{id} :

$$\mathbf{v}_{\text{id}} = \mathcal{D}_I(\mathbf{x}_{\text{src}}), \quad (5)$$

Then, we embed source identity \mathbf{v}_{id} into the diffusion model $\epsilon_{\theta}(\mathbf{x}_t, t, \mathbf{v}_{\text{id}})$, where \mathbf{x}_t is a noisy version of source image \mathbf{x}_{src} at timestep t using Eq. 1.

Loss Functions. Our DiffFace learns a reverse process that reverts the forward process (Fig. 2). Given a noisy source image \mathbf{x}_t , our ID Conditional DDPM aims to predict the noise preserving the source identity \mathbf{v}_{id} by using the denoising score matching loss:

$$\mathcal{L}_{\text{noise}} = \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t, \mathbf{v}_{\text{id}})\|_2^2, \quad (6)$$

where ϵ is a noise added to \mathbf{x}_t . At the same time, we propose an identity loss for the diffusion model to preserve the facial

identity effectively. Since expert models [10, 36, 51, 55] are trained on clean images, we estimate the fully denoised image $\hat{\mathbf{x}}_0$ from each transition \mathbf{x}_t during the denoising diffusion process:

$$\hat{\mathbf{x}}_0 = f_{\theta}(\mathbf{x}_t, t, \mathbf{v}_{\text{id}}), \quad (7)$$

where $f_{\theta}(\cdot)$ is a function to estimate a fully denoised image, which is defined in Eq. 3. Concretely, we induce the identity of predicted denoised image $\hat{\mathbf{v}}_{\text{id}}$:

$$\hat{\mathbf{v}}_{\text{id}} = \mathcal{D}_I(\hat{\mathbf{x}}_0), \quad (8)$$

to be equal to the source identity \mathbf{v}_{id} at every timestep t . Using this predicted identity, the proposed identity loss is as follows:

$$\mathcal{L}_{\text{id}} = 1 - \cos(\mathbf{v}_{\text{id}}, \hat{\mathbf{v}}_{\text{id}}), \quad (9)$$

where $\cos(\cdot, \cdot)$ denotes cosine similarity. Finally, our total loss for ID Conditional DDPM is as follows:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{id}} + \lambda \mathcal{L}_{\text{noise}}. \quad (10)$$

The overall process for training ID Conditional DDPM is summarized in Alg. 1. Also, to verify the benefit of the ID Conditional DDPM, we conduct an ablation study, summarized in Table 2. More examples of our ID Conditional DDPM are in the appendix.

4.2. Facial Guidance

To control the facial attributes of generated images, we propose facial guidance that is applied during the diffusion process. One major advantage of using the diffusion model is that once the model is trained, it can control the image driven by the guidance during the sampling process. Thus we can obtain desired images without any re-training of

Algorithm 2 Diffusion-based Face Swapping with Facial Guidance

Input: source image \mathbf{x}_{src} , target image \mathbf{x}_{targ} , target binary mask M , and masking threshold \hat{T}
Output: swapped image \mathbf{x}_0 that reflects source identity while preserving id-irrelevant attributes in the target image \mathbf{x}_{targ}
 $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $\mathbf{v}_{\text{id}} \leftarrow \mathcal{D}_I(\mathbf{x}_{\text{src}})$
for all t from T to 1 **do**
 $\mu, \sigma \leftarrow \mu_\theta(\mathbf{x}_t), \sigma_\theta(\mathbf{x}_t)$
 $\hat{\mathbf{x}}_0 \leftarrow \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_{\text{id}})}{\sqrt{\bar{\alpha}_t}}$
 $\mathcal{G}_{\text{id}} \leftarrow (1 - \cos(\mathcal{D}_I(\mathbf{x}_{\text{src}}), \mathcal{D}_I(\hat{\mathbf{x}}_0)))$
 $\mathcal{G}_{\text{sem}} \leftarrow \|\mathcal{D}_F(\mathbf{x}_{\text{targ}}) - \mathcal{D}_F(\hat{\mathbf{x}}_0)\|_2^2$
 $\mathcal{G}_{\text{gaze}} \leftarrow \|\mathcal{D}_G(\mathbf{x}_{\text{targ}}) - \mathcal{D}_G(\hat{\mathbf{x}}_0)\|_2^2$
 $\mathcal{G}_{\text{facial}} \leftarrow \lambda_{\text{id}} \mathcal{G}_{\text{id}} + \lambda_{\text{sem}} \mathcal{G}_{\text{sem}} + \lambda_{\text{gaze}} \mathcal{G}_{\text{gaze}}$
 $\hat{\mathbf{x}}_{t-1} \sim \mathcal{N}(\mu - \sigma \nabla_{\mathbf{x}_t} \mathcal{G}_{\text{facial}}, \sigma)$
 $\mathbf{x}_{t-1, \text{targ}} \sim \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_{\text{targ}}, (1 - \bar{\alpha}_{t-1}) \mathbf{I})$
 $M_t \leftarrow \min\{1, \frac{T-t}{T} M\}$
 $\mathbf{x}_{t-1} \leftarrow \hat{\mathbf{x}}_{t-1} \odot M_t + \mathbf{x}_{t-1, \text{targ}} \odot (1 - M_t)$
end for
 $\mathbf{x}_{\text{swap}} \leftarrow \mathbf{x}_0$
return \mathbf{x}_{swap}

the diffusion model. In order to utilize this advantage, we give facial guidance using external models, such as identity embedder [10, 51], face parser [55], and gaze estimator [36] during the sampling process. Note that we can use any off-the-shelf facial model [7, 17], and they can be adaptively selected according to the user’s purpose.

Identity Guidance. In order to transmit the identity information of the source image, we employ the ID Conditional DDPM. However, this identity conditioning turned out to be insufficient, as the synthesized result image kept losing the source’s identity. This is because when only we use the ID Conditional DDPM, the model focuses on other structural information, such as segmentation and gaze, so it fails to maintain source identity. Thus we give identity guidance to prevent loss of source identity in the denoising process. Specifically, we constrained the id vectors of the source and that of generated images to be located closer in the identity embedding space. Our facial identity guidance is formulated as follows:

$$\mathcal{G}_{\text{id}} = 1 - \cos(\mathcal{D}_I(\mathbf{x}_{\text{src}}), \mathcal{D}_I(\hat{\mathbf{x}}_0)). \quad (11)$$

Semantic Guidance. To explicitly match the facial features of the synthesized image to that of the target, we apply the face parsing model [55], which predicts pixel-wise labels for facial components (e.g., nose, eyebrows, and eyes). By leveraging a structural expert, we ensure that the generated image follows the facial structure of the target image.

As shown in Fig. 1, our DiffFace produces a frontal face guided by an ID expert at the beginning of the sampling process. Afterward, our semantic guidance induces the generated image to gradually follow the facial component of

the target image. Finally, the generated image has a similar expression, pose, and shape as the target image by the provided semantic control. In this process, we composed a facial parsing map with essential labels like skin, eyes, and eyebrows, excluding non-facial components such as hair or glasses. We then compute the distance between two selected features. Our semantic guidance is formalized as follows:

$$\mathcal{G}_{\text{sem}} = \|\mathcal{D}_F(\mathbf{x}_{\text{targ}}) - \mathcal{D}_F(\hat{\mathbf{x}}_0)\|_2^2. \quad (12)$$

Gaze Guidance. Gaze information plays a big role in conveying context and emotion. Considering that the face swap technique is heavily used in the film industry, preserving the gaze of the target is a critical issue. Nonetheless, previous models occasionally fail to preserve the target’s gaze. This is because, without explicit gaze modeling, other terms are insufficient to guide gaze. Also, due to the diffusion stochastic process, synthesized results our model tend to have different gazes even if the same input, which makes it impossible to preserve the target’s gaze. Thus, to solve these problems, we explicitly give guidance penalizing different gazes from the target image, using pretrained gaze estimating [36] module.

We first use the off-the-shelf facial landmark-detecting tool to obtain the coordinates of both eyes. Then we crop the eyes of the synthesized image and target image using the coordinates from the previous step. Next, we feed these crop eyes into a pretrained gaze estimating network [36] to obtain the gaze vectors. Finally, we calculate the distance between gaze vectors from the target image and the synthesized image and give it as gaze guidance. Our gaze guidance is formulated as follows:

$$\mathcal{G}_{\text{gaze}} = \|\mathcal{D}_G(\mathbf{x}_{\text{targ}}) - \mathcal{D}_G(\hat{\mathbf{x}}_0)\|_2^2. \quad (13)$$

Incorporation Guidances. To guide the diffusion sampling process toward desired images, we incorporate the gradients from facial guidance modules. As shown in Eq. 4, we can induce facial expert models to behave like classifiers. In particular, given a diffusion model $\epsilon_\theta(\cdot)$ and pretrained facial expert models (e.g., \mathcal{D}_I , \mathcal{D}_F and \mathcal{D}_G), we can derive conditional sampling processes using these facial expert models. Our complete sampling procedure with the incorporated guidance is formulated as follows:

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\mu - \sigma \nabla_{\mathbf{x}_t} \mathcal{G}_{\text{facial}}, \sigma), \quad (14)$$

where

$$\begin{aligned} \mu &= \mu_\theta(\mathbf{x}_t, t, \mathbf{v}_{\text{id}}), & \sigma &= \sigma_\theta(\mathbf{x}_t, t, \mathbf{v}_{\text{id}}), \\ \mathcal{G}_{\text{facial}} &= \lambda_{\text{id}} \mathcal{G}_{\text{id}} + \lambda_{\text{sem}} \mathcal{G}_{\text{sem}} + \lambda_{\text{gaze}} \mathcal{G}_{\text{gaze}}, \end{aligned} \quad (15)$$

and $\mu_\theta(\cdot)$ is the predicted mean of $p_\theta(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{v}_{\text{id}})$. The overall process of facial guidance is summarized in Alg. 2. Note that any additional facial expert models such as 3DMM [4] and face pose estimation [45] can be added to our sampling procedure, which can further enhance the result to preserve the target attributes.

4.3. Target-Preserving Blending

Naïvely using diffusion model in face swap task fails to preserve target background because every region is perturbed by noise. To preserve a background of the target with the desired identity, we utilize a face parser [56] which gives a semantic facial mask. We element-wise product the synthesized result with an obtained target hard mask M to preserve the target’s background while transferring the source facial identity. Similarly, previous methods [3, 29, 30] use a strategy of blending two noisy images with the user-specified binary mask. However, when blending the image in the face swap task with a hard mask, the structural attributes of the target image are removed by noise due to the strong mask intensity.

To achieve our goals, we propose a target-preserving blending method that alters the mask intensity to better preserve structural attributes of target. Target-preserving blending is to gradually increase the mask intensity from zero to one, according to the time of the diffusion process T . By adjusting the starting point where the intensity of the mask becomes one, we can adaptively maintain the structure of the target image. In this manner, we effectively control the trade-off relationship between identity and structure attributes in the face swap task, which will be shown in Sec. 5.4.

Fig. 3 shows an overview of target-preserving blending. M denotes the hard mask obtained by the face parser, and M_t denotes the soft mask whose intensity is strengthened over time. \hat{T} denotes the starting point where the mask intensity becomes one as a hard mask. The soft mask M_t is obtained:

$$M_t = \min(1, \frac{T-t}{\hat{T}}M), \quad (16)$$

We blend the intermediate prediction in the reverse process and target image using the mask M_t . Specifically, we re-define \mathbf{x}_{t-1} in Eq. 14 to $\hat{\mathbf{x}}_{t-1}$ as follows:

$$\hat{\mathbf{x}}_{t-1} \sim \mathcal{N}(\mu - \sigma \nabla_{\mathbf{x}_t} \mathcal{G}_{\text{facial}}, \sigma). \quad (17)$$

Then, to match the noise level between the target image \mathbf{x}_{targ} and the intermediate prediction $\hat{\mathbf{x}}_{t-1}$, we derive a noisy target image $\mathbf{x}_{t-1, \text{targ}}$. Finally, we blend the intermediate prediction $\hat{\mathbf{x}}_{t-1}$ and the noisy target image $\mathbf{x}_{t-1, \text{targ}}$ as

$$\mathbf{x}_{t-1} = \hat{\mathbf{x}}_{t-1} \odot M_t + \mathbf{x}_{t-1, \text{targ}} \odot (1 - M_t). \quad (18)$$

We define \mathbf{x}_{swap} the result of Eq. 18 when $t = 1$. Please refer to Alg. 2 for more details.

5. Experiment

5.1. Implementation Details

Training. Our ID Conditional DDPM follows the architecture of U-Net [43] based on Wide ResNet [57]. We train our ID Conditional DDPM on FFHQ [21] dataset, which

consists of 70k aligned face images with resolution 256×256 . We set $\lambda = 0.5$ at training. We implement our network using PyTorch [37], and AdamW optimizer [28] is used for training, where the learning rate is set to 0.0001. We train the model 700k steps with batch size 48 for approximately 10 days on 8 NVIDIA A100 PCIe 80GB GPUs.

Sampling. In the sampling process, we use extending augmentations used in [3] to prevent adversarial results. The number of extending augmentation is set to 8. We chose $T = 75$ for the number of diffusion steps so that the model could sufficiently alter the target image. We set weights in facial guidance λ_{id} , λ_{sem} , and λ_{gaze} to 2000, 150, and 200, respectively.

5.2. Evaluation Protocol

Evaluation Dataset. We evaluate DiffFace on FaceForensics++ (FF++) dataset [44]. FaceForensics++ dataset is a standard dataset for evaluating face swap methods. All frames are uniformly sampled from the 1000 original videos and generated from various face swap models.

Compared Models. We compare our model with state-of-the-art face swap methods. Because this model is the first to adapt diffusion models to face swap tasks, we evaluate our model with traditional GAN-based methods. We use SimSwap [8], HifiFace [52], InfoSwap [12], MegaFS [60], FaceShifter [25], and Deepfakes [1].

Quantitative Evaluations. We only sample a single image from our model for objective comparison, despite the fact that our model can output diverse images. We quantitatively evaluate models using ID cosine similarity, expression, pose, and shape. For identity metric, we use ArcFace and CosFace embedder to compute the embedding distance between \mathbf{x}_{swap} and \mathbf{x}_{src} . For a fair comparison, we use ArcFace [10] and CosFace [51] to measure models trained with CosFace and ArcFace, respectively. Also, we apply the relative distances metric which was proposed in SmoothSwap [23] in order to measure not only how close \mathbf{x}_{swap} and \mathbf{x}_{src} is, but also how far \mathbf{x}_{swap} and \mathbf{x}_{targ} is. The relative distance metric is formalized in Eq. 19. To evaluate the expression, pose and shape, we use pretrained network [48] to obtain coefficients of 3D face model [26], then compute the distance as follows:

$$D - \mathbf{R} := \frac{\mathcal{D}(\mathbf{x}_{\text{swap}}, \mathbf{x}_{\text{src}})}{\mathcal{D}(\mathbf{x}_{\text{swap}}, \mathbf{x}_{\text{src}}) + \mathcal{D}(\mathbf{x}_{\text{swap}}, \mathbf{x}_{\text{targ}})}, \quad (19)$$

where \mathcal{D} can be any distance metric, e.g. cosine distance.

5.3. Comparison with Baselines

Fig. 4 shows that our DiffFace outperforms other models in terms of changing identity-related attributes. For example, in the first and fourth rows, we notice our result reflects more



Figure 4. **Qualitative comparison of face swap results with other various models.** The results of our model better reflects the source identity while successfully removing target identity-related attributes. Refer Sec. 5.3 for more analysis.

Model	Arc \uparrow	Arc-R \uparrow	Cos \uparrow	Cos-R \uparrow	Expr \downarrow	Pose \downarrow	Shp \downarrow
SimSwap [8]	\dagger	\dagger	0.597	0.756	0.033	0.0005	0.0256
HifiFace [52]	0.575	0.816	0.565	0.792	0.048	0.0007	0.0299
InfoSwap [12]	\dagger	\dagger	0.570	0.841	0.052	0.0010	0.0360
MegaFS [60]	\dagger	\dagger	0.343	0.553	0.046	0.0024	0.0299
FaceShifter [24]	\dagger	\dagger	0.534	0.657	0.061	0.0013	0.0235
DeepFakes [1]	0.443	0.686	0.437	0.635	0.078	0.0022	0.0314
(Cos) DiffFace($\hat{T} = 40$)	0.620	0.859	\dagger	\dagger	0.044	0.0009	0.0269
(Arc) DiffFace($\hat{T} = 40$)	\dagger	\dagger	0.602	0.816	0.043	0.0008	0.0283

Table 1. **Quantitative comparison on FaceForensics++ [44] dataset.** (Arc), (Cos) denotes that model was trained using ArcFace, CosFace respectively. \dagger denotes that model cannot be evaluated because they are trained and evaluated using same identity embedder. Expr, Pose and Shp denotes expression, pose, and shape distance obtained by 3D face model [48] respectively. Please refer to Sec. 5 for more details.

vivid lips and eyes, while other results models tend to have eyes and lip colors from target images. This shows that our model more effectively transfers identity-related attributes than other models.

Table 1 shows the same tendency as the qualitative result shown above. As shown in the first four columns, our model achieves the highest identity score. Also, the ability to remove the target’s identity is superior to any other model. We speculate that these results are caused by the unique property of the diffusion process. Specifically, existing GAN-based models could not remove specific regions, while diffusion

Model	Arc \uparrow	Arc-R \uparrow	Cos \uparrow	Cos-R \uparrow	Expr \downarrow	Pose \downarrow	Shp \downarrow
(I) Unconditional	0.486	0.803	0.455	0.801	0.063	0.0016	0.0358
(II) Conditional + ($\lambda_{id} = 0$)	0.455	0.751	0.535	0.781	0.040	0.0008	0.0220
(III) (Cos) $\hat{T} = 30$	0.598	0.813	\dagger	\dagger	0.037	0.0007	0.0226
(IV) (Cos) $\hat{T} = 40$	0.620	0.859	\dagger	\dagger	0.044	0.0009	0.0269
(V) (Cos) $\hat{T} = 50$	0.634	0.888	\dagger	\dagger	0.050	0.0011	0.0303
(VI) (Arc) $\hat{T} = 30$	\dagger	\dagger	0.580	0.766	0.035	0.0006	0.0240
(VII) (Arc) $\hat{T} = 40$	\dagger	\dagger	0.602	0.816	0.043	0.0008	0.0283
(VIII) (Arc) $\hat{T} = 50$	\dagger	\dagger	0.603	0.816	0.049	0.0009	0.0311

Table 2. **Quantitative ablation study.** Performance of various configurations on FF++. $\lambda_{id} = 0$ denotes that identity guidance was not used, and \hat{T} denotes the timestep where mask intensity becomes one. For other notations, please refer the caption of Table 1.

models can remove target information with random noise.

Although our approach yields limited performance on expression, pose, and shape score, one thing to note here is that we can manipulate guidance in the sampling process in order to compensate the problem. Specifically, we can adjust the reflectivity of the target image in the reverse process by controlling \hat{T} , which is the starting point where the intensity of the mask becomes one. In summary, we can dynamically alter our model to satisfy different purposes. Please refer to Sec. 5.4 for more details.

5.4. Ablation Study and Analysis

Identity Guidance. (I) and (II) in Table 2 show a quantitative ablation study on ID Conditional DDPM and identity guidance. Our model results in the highest identity similarity

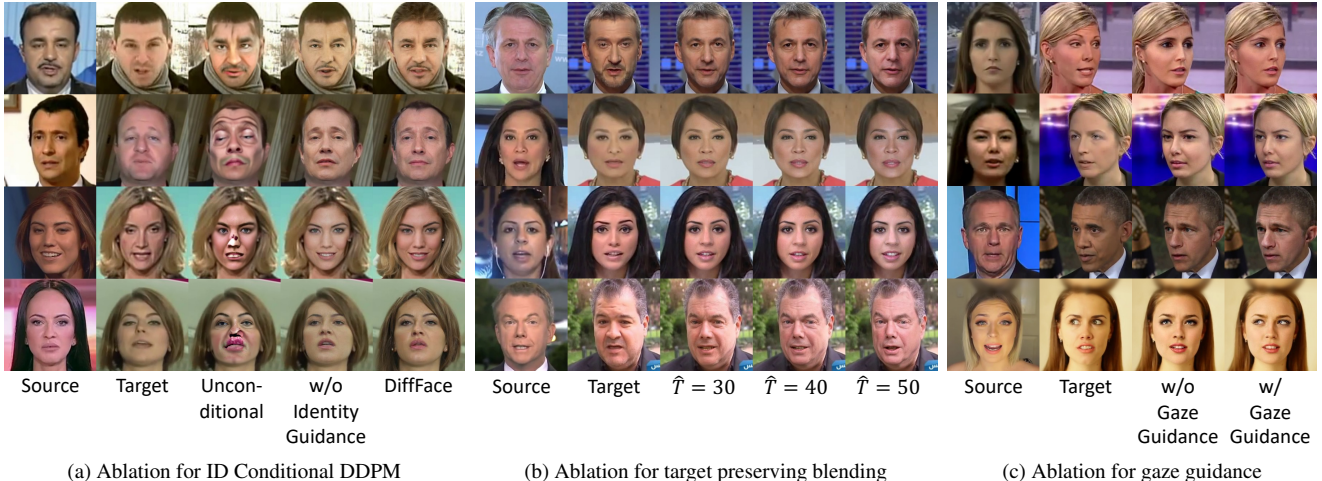


Figure 5. **Qualitative results of ablations:** (a) ID Conditional DDPM, (b) target preserving blending, and (c) gaze guidance. Our method can be easily controlled without any additional training. We can obtain various face swapped results by controlling the facial module.

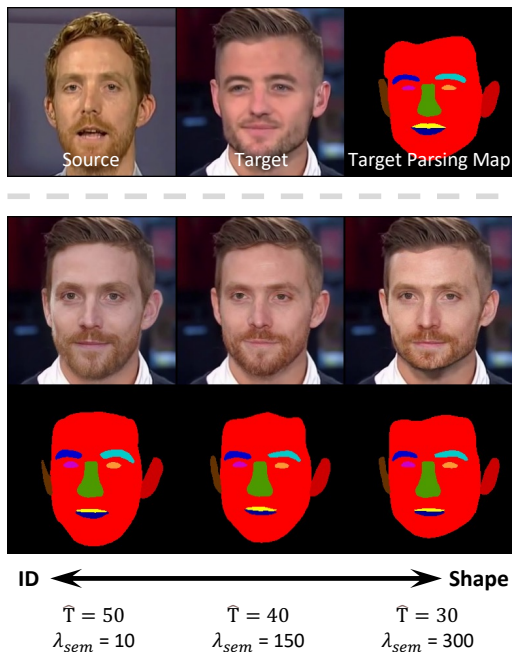


Figure 6. **Trade-off between ID and shape.** The results show how our model can be altered to prioritize either shape or id. The second and third row represents the synthesized image and corresponding face parsing map.

score when ID Conditional DDPM and identity guidance are present. The same trend can be found in Fig. 5 (a), which shows a qualitative ablation study on ID Conditional DDPM and identity guidance. It is noticeable that results without ID Conditional DDPM show visible artifacts, and results without identity guidance show poor identity transfer compared to the full model.

Target Preserving Blending. In this ablation study, we adjust \hat{T} , the starting point where the mask becomes one

as a hard mask, to investigate the effect of target preserving blending, as shown in (III) to (VIII) of Table 2. As \hat{T} increases, the id similarity score also increases because the model gains more control over the image and thus manipulates the image close to the source image. At the same time, expression, pose, and shape deteriorate because the image gets manipulated more than before. The same tendency is also shown in Fig. 5 (b). As \hat{T} decreases, the synthesized image tends to move toward the target image, in terms of expression, pose and shape. Also, we can observe that the skin tone of the synthesized image approaches the target. In summary, we can choose our model to balance between preserving the target attributes and injecting source identity without additional training.

Gaze Guidance. Fig. 5 (c) illustrates our qualitative ablation study on gaze guidance. Recall that the synthesized image’s gaze must follow the target image. We can find that with gaze guidance, the synthesized image has the same gaze as the target image. Thus we employ gaze guidance to fulfill gaze information which was not controllable using the identity or semantic guidance.

Trade-off Control on ID and Shape. Fig. 6 shows how our model can be altered to prioritize either ID or shape. We can induce our model to show active shape change by incrementing \hat{T} and decrementing λ_{sem} . Conversely, we can force the model to focus on preserving target structure by decrementing \hat{T} and incrementing λ_{sem} . As a result, we can see that the shape of the chin on the left is gradually getting rounder as it goes to the right. In this way, we can easily have controllability without the need to retrain the model in the face swap task.

6. Conclusion

In this work, we propose *DiffFace*, a diffusion-based framework aiming for controllable and high-fidelity subject-agnostic face swapping. Based on the trained our ID Conditional DDPM, the facial guidance from the pretrained facial expert models make the model to faithfully transfer the source identity while preserving target attributes during the sampling process. The target preserving blending helps our framework to keep the attributes of the target face from noise while transferring the source facial identity. Moreover, our framework can flexibly apply various facial guidance and adaptively control the ID-attribute tradeoff to achieve better results without additional training. Extensive experiments demonstrate that the proposed framework significantly outperforms previous face swapping methods.

References

- [1] DeepFakes (<https://github.com/deepfakes/faceswap>), Nov. 2021. [2](#), [6](#), [7](#)
- [2] Xiang An, Jiangkang Deng, Jia Guo, Ziyong Feng, Xuhan Zhu, Yang Jing, and Liu Tongliang. Killing two birds with one stone: Efficient and robust training of face recognition cnns by partial fc. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. [12](#), [14](#)
- [3] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. [2](#), [6](#)
- [4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '99*, pages 187–194, USA, July 1999. ACM Press/Addison-Wesley Publishing Co. [2](#), [5](#)
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2019. [2](#)
- [6] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer Vision*, 2017. [12](#)
- [7] Feng-Ju Chang, Anh Tran, Tal Hassner, Iacopo Masi, Ram Nevatia, and Gérard Medioni. Expnet: Landmark-free, deep, 3d facial expressions. In *13th IEEE Conference on Automatic Face and Gesture Recognition*, 2018. [5](#)
- [8] Renwang Chen, Xuanhong Chen, Bingbing Ni, and Yanhao Ge. Simswap: An efficient framework for high fidelity face swapping. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2003–2011, 2020. [2](#), [3](#), [6](#), [7](#)
- [9] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. *arXiv preprint arXiv:2108.02938*, 2021. [2](#)
- [10] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019. [4](#), [5](#), [6](#), [12](#), [14](#)
- [11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. [2](#), [12](#)
- [12] Gege Gao, Huaibo Huang, Chaoyou Fu, Zhaoyang Li, and Ran He. Information bottleneck disentanglement for identity swapping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3404–3413, 2021. [2](#), [3](#), [6](#), [7](#)
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. [2](#)
- [14] 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. [12](#)
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. [2](#), [3](#)
- [16] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. [2](#)
- [17] Feng ju Chang, Anh Tran, Tal Hassner, Iacopo Masi, Ram Nevatia, and Gérard Medioni. FacePoseNet: Making a case for landmark-free face alignment. In *7th IEEE International Workshop on Analysis and Modeling of Faces and Gestures, ICCV Workshops*, 2017. [5](#)
- [18] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation, 2017. [2](#)
- [19] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *Advances in Neural Information Processing Systems*, 33:12104–12114, 2020. [13](#)
- [20] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. [3](#)
- [21] 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. [2](#), [3](#), [6](#)
- [22] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. [3](#)
- [23] Jiseob Kim, Jihoon Lee, and Byoung-Tak Zhang. Smooth-Swap: A Simple Enhancement for Face-Swapping with Smoothness. *arXiv:2112.05907 [cs]*, Dec. 2021. [6](#)
- [24] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457*, 2019. [2](#), [7](#)

- [25] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. FaceShifter: Towards High Fidelity And Occlusion Aware Face Swapping. *arXiv:1912.13457 [cs]*, Dec. 2019. [2](#), [6](#)
- [26] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), 2017. [6](#)
- [27] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu, Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with semantic diffusion guidance. *arXiv preprint arXiv:2112.05744*, 2021. [2](#)
- [28] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2017. [6](#)
- [29] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. [2](#), [6](#)
- [30] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. [2](#), [6](#)
- [31] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *CoRR*, abs/1411.1784, 2014. [2](#)
- [32] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks, 2018. [2](#)
- [33] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. [2](#)
- [34] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. pages 8162–8171, 2021. [2](#)
- [35] 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. [2](#)
- [36] Seonwook Park, Xucong Zhang, Andreas Bulling, and Otmar Hilliges. Learning to find eye region landmarks for remote gaze estimation in unconstrained settings. In *ACM Symposium on Eye Tracking Research and Applications (ETRA)*, ETRA '18, New York, NY, USA, 2018. ACM. [4](#), [5](#), [12](#)
- [37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. [6](#)
- [38] Justin NM Pinkney and Doron Adler. Resolution dependent gan interpolation for controllable image synthesis between domains. *arXiv preprint arXiv:2010.05334*, 2020. [13](#), [18](#)
- [39] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Diffusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10619–10629, 2022. [2](#)
- [40] 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. [2](#)
- [41] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. [2](#), [3](#), [4](#)
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. [2](#), [4](#)
- [43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015. [6](#)
- [44] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Niessner. FaceForensics++: Learning to Detect Manipulated Facial Images. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1–11, Seoul, Korea (South), Oct. 2019. IEEE. [6](#), [7](#), [12](#), [13](#), [16](#), [17](#)
- [45] Nataniel Ruiz, Eunji Chong, and James M. Rehg. Fine-grained head pose estimation without keypoints. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018. [5](#)
- [46] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. [2](#), [4](#)
- [47] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. [2](#), [4](#)
- [48] Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael Black. Learning to regress 3d face shape and expression from an image without 3d supervision. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [6](#), [7](#)
- [49] Junyoung Seo, Gyuseong Lee, Seokju Cho, Jiyoung Lee, and Seungryong Kim. Midms: Matching interleaved diffusion models for exemplar-based image translation. *arXiv preprint arXiv:2209.11047*, 2022. [2](#)
- [50] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. [2](#)
- [51] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *Proceedings of*

- the IEEE conference on computer vision and pattern recognition*, pages 5265–5274, 2018. 4, 5, 6
- [52] Yuhan Wang, Xu Chen, Junwei Zhu, Wenqing Chu, Ying Tai, Chengjie Wang, Jilin Li, Yongjian Wu, Feiyue Huang, and Rongrong Ji. Hiface: 3d shape and semantic prior guided high fidelity face swapping. *arXiv preprint arXiv:2106.09965*, 2021. 2, 6, 7
- [53] Yangyang Xu, Bailin Deng, Junle Wang, Yanqing Jing, Jia Pan, and Shengfeng He. High-resolution face swapping via latent semantics disentanglement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7642–7651, 2022. 3
- [54] Zhiliang Xu, Hang Zhou, Zhibin Hong, Ziwei Liu, Jiaming Liu, Zhizhi Guo, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Styleswap: Style-based generator empowers robust face swapping. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 3
- [55] Bisenet Yu, J Wang, C Peng, C Gao, G Yu, N Sang, and Bisenet. Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, page 325. 4, 5
- [56] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *European Conference on Computer Vision*, pages 334–349. Springer, 2018. 6, 12
- [57] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 6
- [58] Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations. *arXiv preprint arXiv:2207.06635*, 2022. 2
- [59] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A. Efros. Generative visual manipulation on the natural image manifold. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2016. 2
- [60] Yuhao Zhu, Qi Li, Jian Wang, Cheng-Zhong Xu, and Zhenan Sun. One shot face swapping on megapixels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4834–4844, 2021. 2, 3, 6, 7

Appendix

A. Architecture Details

A.1. Identity Embedding Model

We employ a pretrained ArcFace [10] and CosFace [2] identity embedding models to extract the identity vector of the source face. ArcFace and CosFace models are based on ResNet-101 [14] and the identity vector becomes an unit-length ($\|v_{id}\| = 1$) by the *UnitNorm* in Fig. 7.

A.2. ID Conditional DDPM

Our diffusion model follows the architecture of guided diffusion model [11], except for conditioning the identity. In order to sample results with the desired identity, we inject the identity vector obtained from the identity embedder [2, 10] into the residual blocks of U-Net. As a result, we can condition our diffusion model and sample images with specific identities. The detailed structure of ID Conditional DDPM is depicted in Fig. 7. Also, we show some examples of ID Conditional DDPM in Fig. 9.

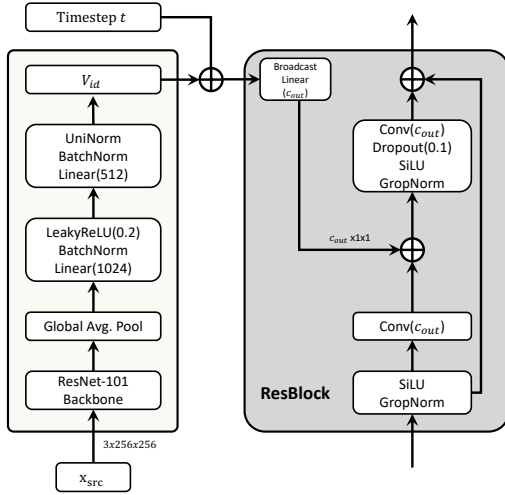


Figure 7. The residual block structure of ID Conditional DDPM. The identity vector v_{id} is embedded into each Resblocks of U-Net with the diffusion timestep t . We denote \oplus the summation operation and feature dimensions are written in the order of (channels \times height \times width)

B. Sampling Details

B.1. Face Parser

In every timestep of the sampling process, we use pretrained face parsing network [56] to give facial guidance. The pretrained face parsing network outputs 19 classes, including all facial components and accessories. We construct a facial mask using 11 classes related to the face swap task (e.g.,

skin, nose, and eyes) and exclude non-facial components (e.g., hair, a hat, and cloth).

B.2. Gaze Estimator

Similar to face parser, we give gaze guidance at every time step during the sampling process. We first use facial landmark detecting tool [6] to obtain coordinates of both eyes and then crop eye images. The eye images are automatically resized to 96x160. Then the images are fed into the pretrained gaze estimator [36] to obtain gaze-relevant features.

C. Additional Results from DiffFace

C.1. Target-Preserving Blending

Fig. 10a shows the change of the noisy image x_t and facial mask M_t during the diffusion process. Our facial mask intensity increases linearly until the masking threshold \hat{T} . By altering the mask’s intensity, which determines the reflectivity of the target image, we can preserve the fine details of the facial attributes (e.g., expression, pose and skin tone).

C.2. Trade-off between ID and Shape

Unlike GAN-based face swapping method that produces a deterministic output image, our method can control the face shape with the facial guidance $\mathcal{G}_{\text{facial}}$ from λ_{sem} and masking threshold \hat{T} . Fig. 10b shows the result of emphasizing ID or face shape. For example, the sample in the first column of the third row shows a round chin shape, while the sample in the fourth column of the third row shows a sharp chin shape. Also, we can observe the skin tone changes in various ways according to the time threshold \hat{T} .

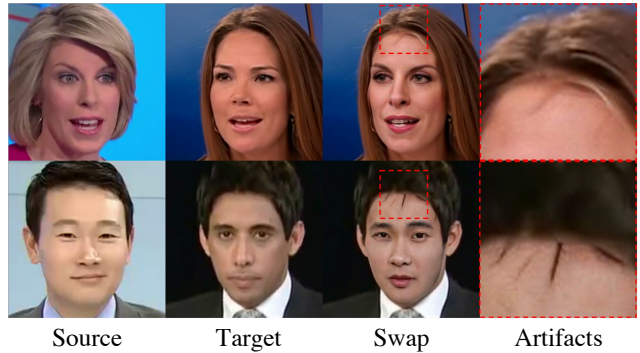


Figure 8. Some failure cases of DiffFace on the FaceForensics++ dataset [44].

C.3. Comparison and Out-Of-Domain Results

We provide additional collections of swapped-image samples based on our DiffFace model. Fig. 11 show the extra face swapping results of various models on the FaceForensic++

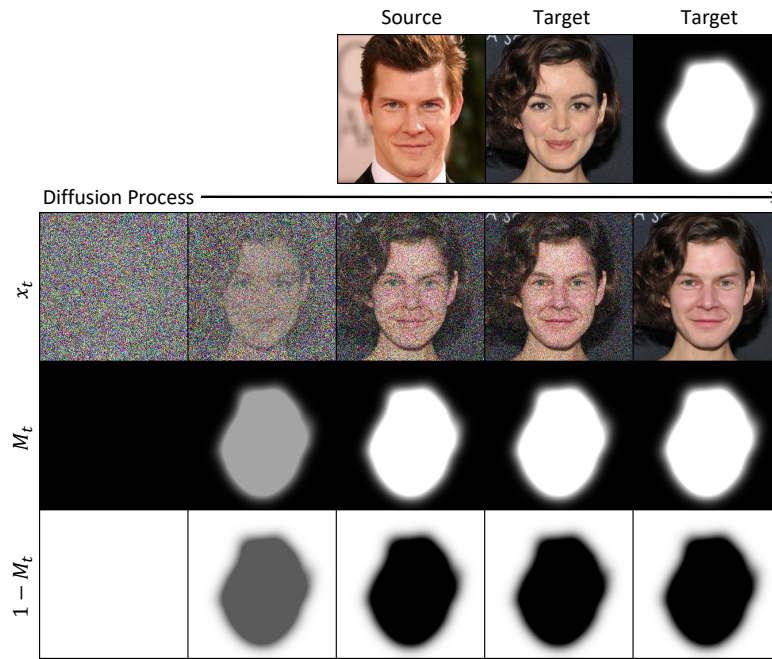
dataset [44]. Fig. 12 and Fig. 13 show the results of out-of-domain datasets, where oil portrait paintings (Metfaces dataset [19]) and cartoon faces (Disney Face dataset [38]) are used. Although our model is not trained on any of these images, the results reflect the characteristics of each domain with shape changing.

D. Limitations

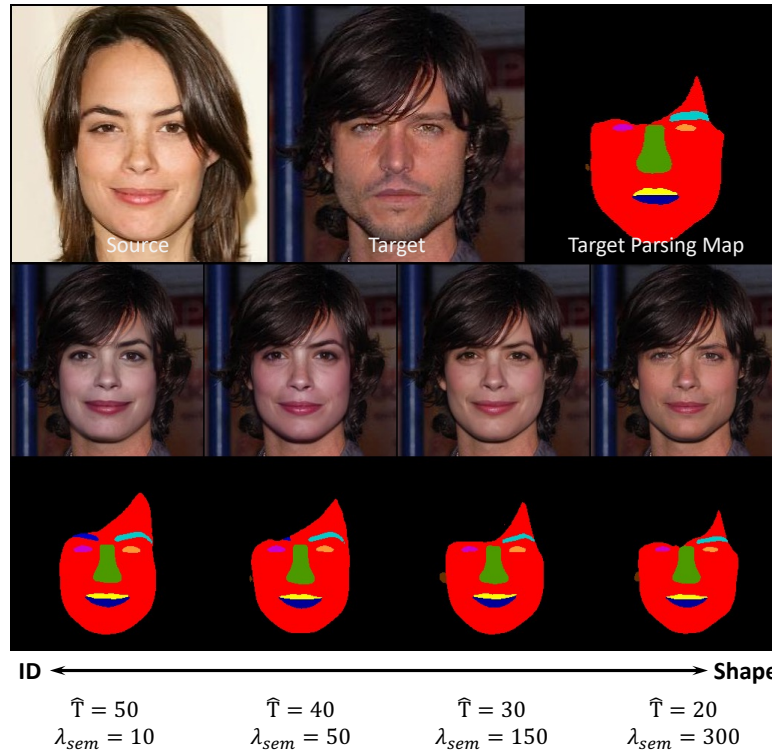
Denosing diffusion probabilistic models have shown remarkable generative performance in various computer vision tasks. However, since the diffusion model generates images through sequential stochastic transitions, once an artifact (i.e., wrinkles, hair segments, and glasses) occurs, it is easy to maintain it (Fig. 8). For this reason, our model sometimes shows unintended artifacts in the output. Therefore, research to correcting transitions with artifacts seems to be necessary.



Figure 9. Image samples from ID Conditional DDPM (Sec. A.2). The first row is the source faces which are provided to the identity embedder [2, 10]. The last five rows are samples generated by ID Conditional DDPM.



(a) Visualization result of a noisy image x_t , facial mask M_t and background mask $1 - M_t$ for the reverse process.



(b) Visualization result for the trade-off between ID and Shape by controlling \hat{T} and λ_{sem}

Figure 10. More ablation results for Target-Preserving Blending and Trade-off between ID and Shape



Figure 11. Additional comparison results of diverse models on the FaceForensics++ dataset [44].

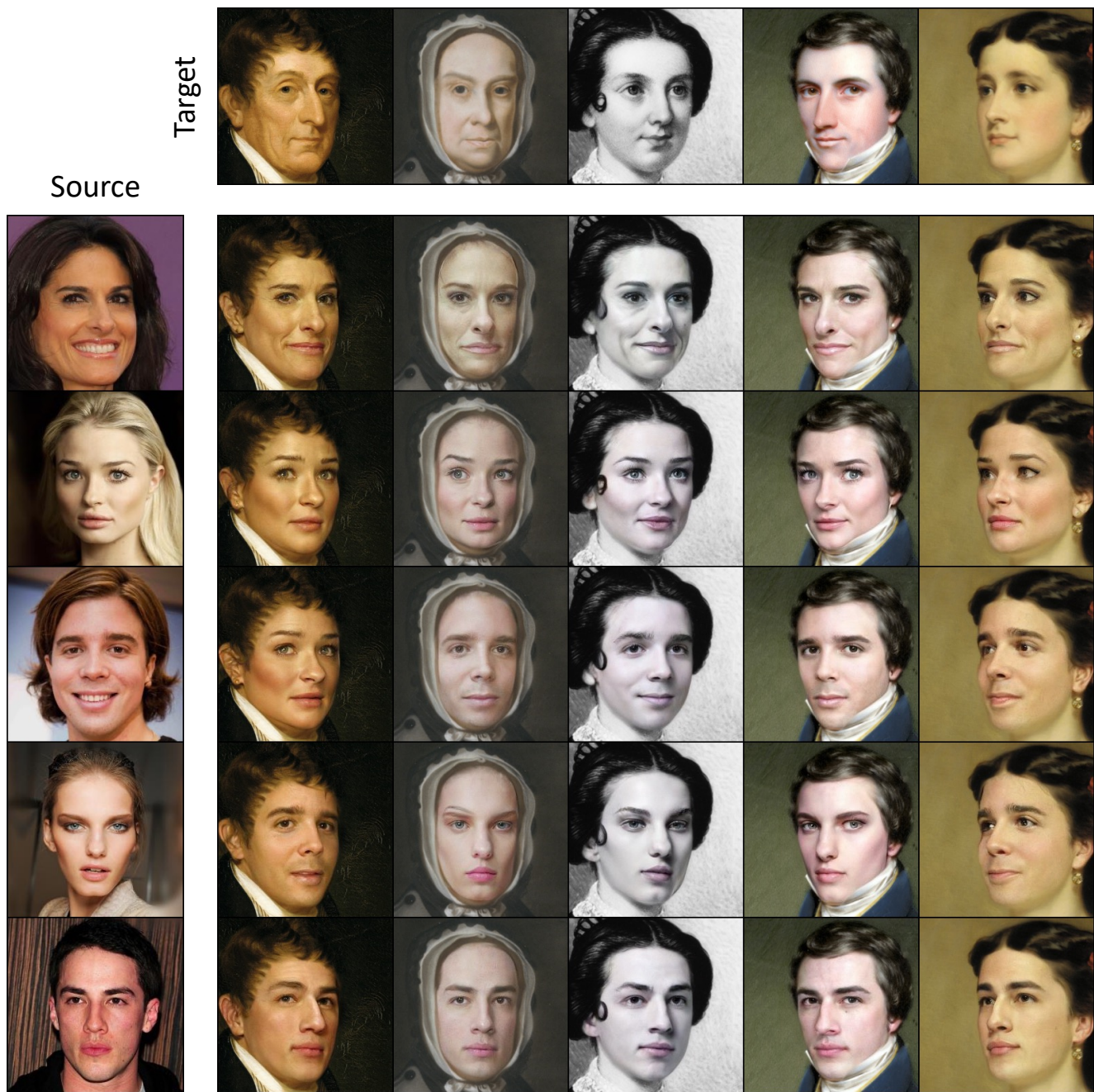


Figure 12. Out-of-domain face swapping results generated by our DiffFace on the MetFaces dataset [44].

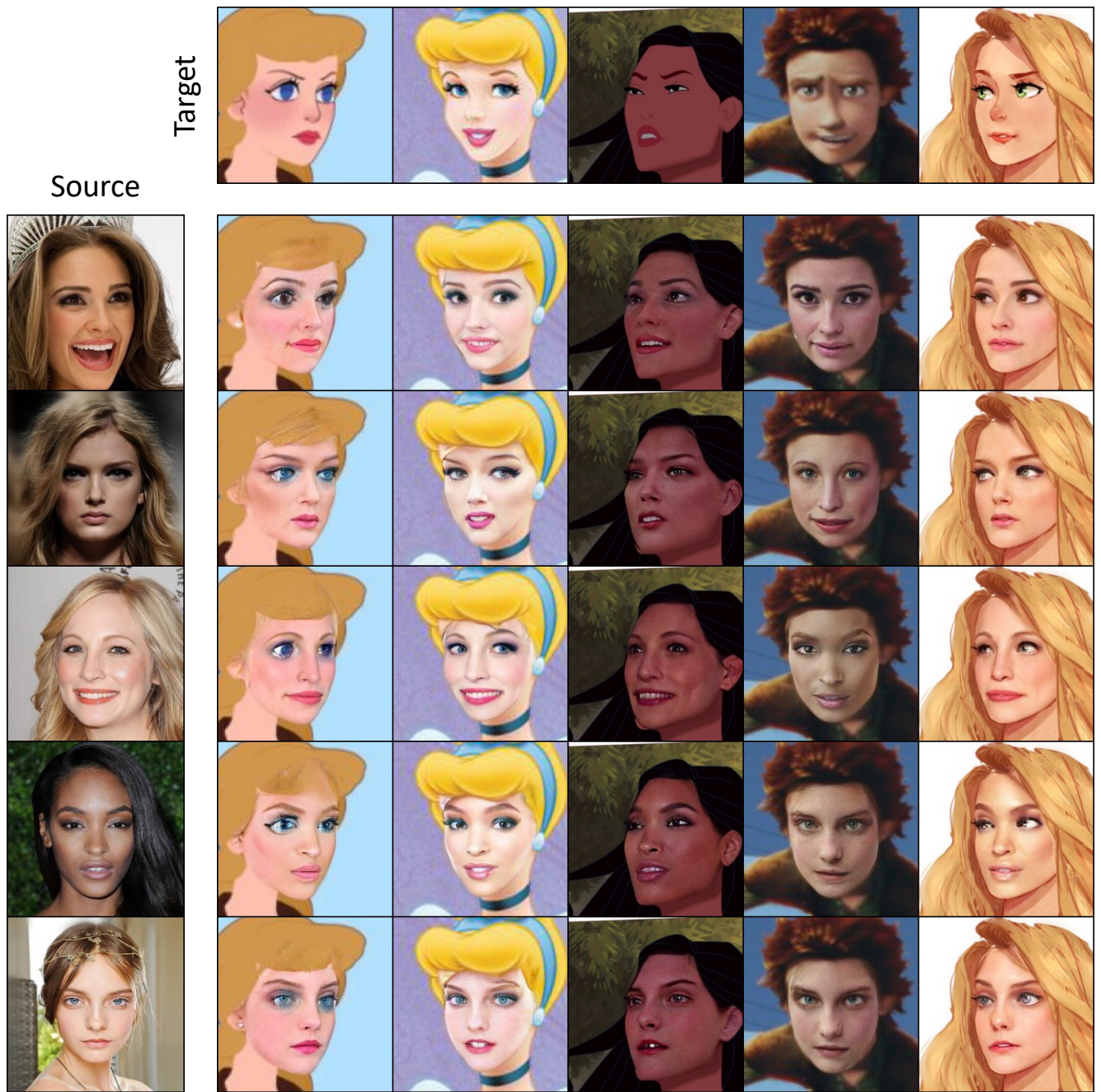


Figure 13. Out-of-domain face swapping results generated by our DiffFace on the Disney Face dataset [38].