

# Cross-Domain and Disentangled Face Manipulation with 3D Guidance

CAN WANG, City University of Hong Kong  
MENGLI CHAI, Snap Inc.  
MINGMING HE, USC Institute for Creative Technologies  
DONGDONG CHEN, Microsoft Cloud AI  
JING LIAO\*, City University of Hong Kong

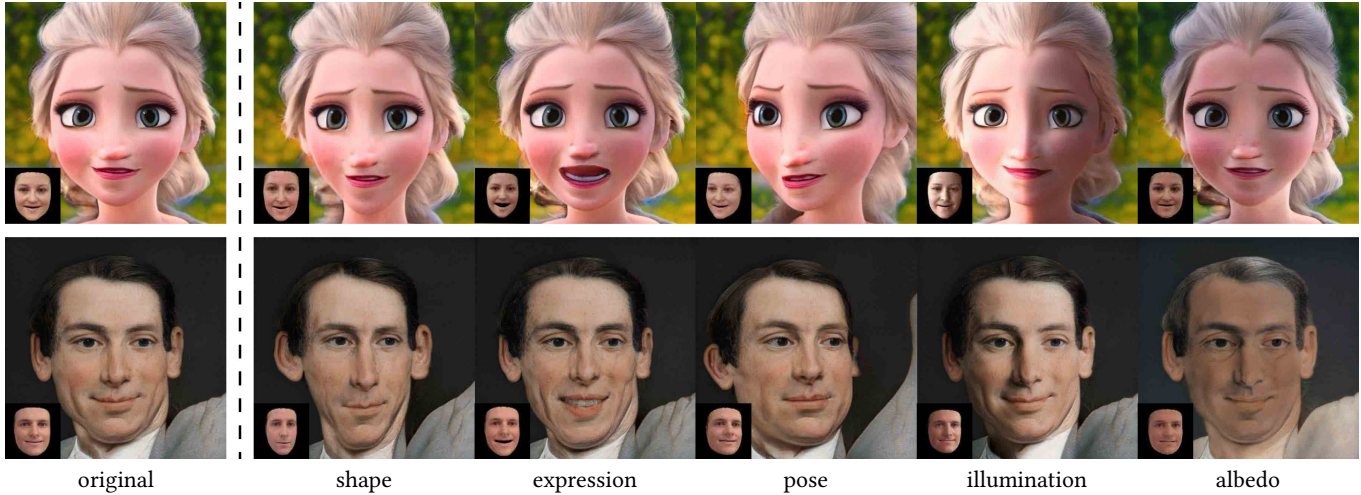


Fig. 1. Our cross-domain face manipulation pipeline allows disentangled manipulation of arbitrary out-of-domain face images. Controlled by 3D human face parameters (show as the inset in each result), our method enables editing of a wide variety of semantic facial attributes, including shape, expression, pose, illumination, and albedo.

Face image manipulation via three-dimensional guidance has been widely applied in various interactive scenarios due to its semantically-meaningful understanding and user-friendly controllability. However, existing 3D-morphable-model-based manipulation methods are not directly applicable to out-of-domain faces, such as non-photorealistic paintings, cartoon portraits, or even animals, mainly due to the formidable difficulties in building the model for each specific face domain. To overcome this challenge, we propose, as far as we know, the first method to manipulate faces in arbitrary domains using human 3DMM. This is achieved through two major steps: 1) disentangled mapping from 3DMM parameters to the latent space embedding of a pre-trained StyleGAN2 [Karras et al. 2020] that guarantees disentangled and precise controls for each semantic attribute; and 2) cross-domain adaptation that bridges domain discrepancies and makes human 3DMM applicable to out-of-domain faces by enforcing a consistent latent space embedding. Experiments and comparisons demonstrate the superiority of our high-quality semantic manipulation method on a variety of face domains with all major 3D facial attributes controllable – pose, expression, shape, albedo, and illumination. Moreover, we develop an intuitive editing interface to support user-friendly control and instant feedback. Our project page is <https://cassiepython.github.io/sigasia/cddfm3d.html>.

CCS Concepts: • **Computing methodologies** → **Image processing**.

\*Corresponding Author.  
Authors' addresses: Can Wang, City University of Hong Kong, [cwang355-c@my.cityu.edu.hk](mailto:cwang355-c@my.cityu.edu.hk); Menglei Chai, Snap Inc. [cmlatsim@gmail.com](mailto:cmlatsim@gmail.com); Mingming He, USC Institute for Creative Technologies, [hmm.lillian@gmail.com](mailto:hmm.lillian@gmail.com); Dongdong Chen, Microsoft Cloud AI, [cddlyf@gmail.com](mailto:cddlyf@gmail.com); Jing Liao\*, City University of Hong Kong, [jingliao@cityu.edu.hk](mailto:jingliao@cityu.edu.hk).

Additional Key Words and Phrases: Face Image Manipulation, Domain Adaptation, 3D Morphable Model, Disentanglement, StyleGAN

## 1 INTRODUCTION

Face image manipulation has been long coveted in various interactive scenarios, such as photo enhancement, social games, and virtual reality. With the striking development of human digitization techniques, 3D-guided face manipulation has popularized itself through semantically-meaningful understanding and user-friendly controllability. With the strong parameterization capability of 3D morphable models (3DMMs) [Blaiz and Vetter 1999], there arises the possibility of editing various facial attributes, including pose, expression, shape, albedo, and illumination. Along this direction, before the deep learning era, the problem is mainly addressed from the perspective of the traditional graphics pipeline, in which a face model is fitted to the subject in the image and then re-rendered with altered facial parameters. However, due to the limitations of the low dimensionality of 3DMM and the approximated shading models, the results often suffer from unsatisfactory realism. More recently, the rapid advancement of deep generative networks, such as StyleGAN [Karras et al. 2019] and StyleGAN2 [Karras et al. 2020], enables photo-realistic face synthesis. To bring the best of both worlds, the pioneering works of StyleRig [Tewari et al. 2020a] and PIE [Tewari et al. 2020b] achieve huge success via adopting 3DMM as the proxy

of intuitive parametric controls and relying on pre-trained StyleGAN to synthesis the corresponding high-quality manipulation results.

Despite the progress made in 3D-guided face manipulation, most existing methods are available only for human faces and cannot be trivially extended to out-of-domain faces, such as non-photorealistic paintings, cartoon portraits, or even animals. This is mainly due to the formidable difficulties in building the 3DMM for each specific face domain, regarding both data acquisition and processing. To overcome this challenge, we propose, as far as we know, the first method to manipulate arbitrary out-of-domain faces via human 3DMM. This is achieved through two major steps: *disentangled attribute-latent mapping* and *cross-domain adaptation*.

Within the human-face domain, the disentangled attribute-latent mapping structure is adopted to bridge between the latent space and the 3DMM parameters. Specifically, a source latent code, inverted from the input image, is first mapped to its 3DMM parameters for manipulation. Then, the edited parameters are projected back to update the latent code that generates the final image. Different from StyleRig [Tewari et al. 2020a] that treats the latent code of StyleGAN2 as an indivisible element in mapping, we propose the *reduced StyleSpace*, which is decomposed into subspaces corresponding to different semantic 3DMM attributes, i.e., pose, expression, shape, albedo, and illumination, respectively. Based on that, the aforementioned mapping is disentangledly learned between each attribute and its corresponding subspace. Our method minimizes the mutual interference between different attributes and thus greatly improves the quality and flexibility of the manipulation.

Extending to out-of-domain faces, our cross-domain adaptation makes the semantic latent embedding consistent for all domains. Inspired by the recent observation in [Huang et al. 2020; Pinkney and Adler 2020], we perform latent-consistent fine-tuning that adapts the StyleGAN2 generator to another domain while keeping the latent spaces aligned. Therefore, given an out-of-domain face image, we first optimize its corresponding latent code through cross-domain latent inversion and then manipulate the code using the in-domain attribute-latent mapping that is consistent for all domains. Furthermore, when fed into the original human-face StyleGAN generator, the inverted latent code can be mapped to a human face. As a manipulation proxy, any edits on it will be faithfully reflected on the out-of-domain face input by simply feeding the same edited latent code into the finetuned StyleGAN2 for out-of-domain faces, thanks to the shared latent embedding.

Our method enables high-quality semantic manipulation on a variety of out-of-domain faces with all major 3D facial attributes controllable – pose, expression, shape, albedo, and illumination, as shown in Fig. 1. Empowered by StyleGAN2, these edits maintain visually-realistic occlusion handling, lighting consistency, and perspective parallax, which can be hardly achieved with traditional approaches. Furthermore, with our reduced latent space and disentangled attribute-latent mapping, our method enjoys highly-disentangled manipulation, which allows it to edit one attribute without affecting irrelevant content or manipulate multiple attributes simultaneously without introducing mutual interference. Experiments and comparisons demonstrate the superiority of the proposed method compared to previous manipulation methods (either in image space [Schaefer

et al. 2006; Siarohin et al. 2019b] or latent space [Härkönen et al. 2020; Shen et al. 2020]), regarding both quality and controllability. In addition, we develop an editing interface for user-friendly manipulation based on the 3D face proxy, which achieves intuitive controls and instant feedback.

The contributions of this paper can be summarized as follows.

- We propose the first method to manipulate semantic attributes of out-of-domain faces with human face 3DMM;
- Our reduced StyleSpace and disentangled attribute-latent mapping guarantee disentangled and precise controls of all facial attributes, achieving disentangled and high-quality face manipulation;
- Our cross-domain adaptation approach bridges domain discrepancies and makes human face 3DMM applicable to out-of-domain faces by enforcing a consistent latent space embedding.

Our source code, pre-trained models, and data will be made publicly available to facilitate future research.

## 2 RELATED WORK

### 2.1 Image Space Face Manipulation

Before the deep learning era, many traditional methods have been proposed for face manipulation in the image space. For example, Moving Least Squares (MLS) [Schaefer et al. 2006] uses sparse control points or lines to create smooth deformations of images. The warping-based methods [Averbuch-Elor et al. 2017; Fried et al. 2016] approximate the 3D space editing effects via lightweight 2D warps in the image plane. Recently, deep neural networks have been extensively exploited towards warping-based facial geometry and expression synthesis. For instance, Geng et al. [2018] use a GAN to synthesize appearance details onto a pre-warped face image. Wiles et al. [2018] propose to generate a dense motion field by using neural networks to warp the image towards a reference face. Siarohin et al. [2019a] and Siarohin et al. [2019b] propose to encode motion information via keypoints learned in a self-supervised fashion and then warp the image according to the reference keypoint trajectories.

Besides keypoints, facial landmarks are also popular to represent face geometry information, widely used for conditional face translation [Ha et al. 2020; Natsume et al. 2018; Nirkin et al. 2019; Xiang et al. 2020; Zakharov et al. 2019]. The one-shot subject-agnostic frameworks by Natsume et al. [2018] and Nirkin et al. [2019] can achieve face swapping between unseen identities conditioned on 2D landmarks. And Zakharov et al. [2019] show that few-shot learning can improve the quality of portrait reenactment, but it cannot adapt the landmarks between different subjects and results in identity mismatch. To alleviate this problem, some methods Ha et al. [2020]; Xiang et al. [2020] further introduce landmark disentanglement which isolates the geometry from the identity. However, without the guidance from 3D face geometry, such 2D image editing methods still fail to guarantee consistent quality under large attribute changes, including poses, shapes, expressions, and illuminations.

### 2.2 Latent Space Face Manipulation

Generative adversarial networks (GANs) [Arjovsky et al. 2017; Goodfellow et al. 2014] have achieved great success in recent years. The

most popular GAN models for face generation are the recent work StyleGAN [Karras et al. 2019] and its improved version StyleGAN2 [Karras et al. 2020], which can synthesize high fidelity face images from one learned latent space. However, it is difficult to control the synthesis process in a semantically controllable and disentangled manner, *i.e.*, changing one attribute (*e.g.* pose) while preserving consistency of other attributes (*e.g.* identity, expression, background). To overcome this limitation, some recent works [Härkönen et al. 2020; Shen et al. 2020; Shen and Zhou 2020] explore how to manipulate the underlying learned latent space. Essentially, they are all related to the vector arithmetic property observed in [Radford et al. 2016], *i.e.*, semantic editing operations can be achieved by first computing a difference vector between two latent vectors and then adding it onto another latent vector. In other words, to change one specific attribute, the key is to find the directional vector corresponding to this attribute.

Specifically, InterFaceGAN [Shen et al. 2020] proposes to learn a binary attribute classifier to identify the separation boundary for each attribute and further require the separation boundaries of different attributes to be as orthogonal as possible. Such attribute classifiers are also used in the very recent work StyleFlow [Abdal et al. 2020], which formulates conditional exploration as conditional continuous normalizing flows in the latent space. They share the common limitation that a labeled dataset containing all the attribute labels is required, which also limits their supported attribute types, especially for out-of-domain faces. In contrast, the noteworthy work GANSpace [Härkönen et al. 2020] can identify important latent directions for different attributes in a unsupervised way. It applies principle component analysis either in latent space or feature space, and achieves interpretable control with layer-wise perturbation along the principal directions. To transfer such controllability to real image editing, the only extra step is to map the real image into the latent space via some GAN inversion methods [Abdal et al. 2019; Zhu et al. 2020]. Broadly speaking, our method also belongs to latent space manipulation, but we provide more fine-grained and rig-like control like [Tewari et al. 2020a] by building a bidirectional mapping between the latent space and the 3DMM space.

### 2.3 3D Guided Face Manipulation

3D Morphable Models (3DMM) originally proposed by [Blanz and Vetter 1999] represent human faces in a parametric space. It provides a flexible and explicit control over 3D facial attributes (*e.g.* geometry, expression, and pose), and further enables rendering of reconstructed faces with illumination modeled via spherical harmonics parameters [Ramamoorthi and Hanrahan 2001]. The reconstructed face models can be used in a wide range of face manipulation applications, such as video dubbing [Dale et al. 2011; Garrido et al. 2015], face reenactment [Thies et al. 2015; Zollhöfer et al. 2018], and lip-syncing [Suwajanakorn et al. 2017]. However, due to the use of mostly linear statistical models, the reconstructed faces only capture the coarse geometry shape, with subtle details missing. Also, computer graphics rendering often leads to non-photorealistic results. To enhance the photorealism, neural rendering techniques propose to build differentiable graphics pipelines [Kim et al. 2019, 2018; Tewari et al. 2017], or combine the traditional graphics pipeline

with learnable components [Thies et al. 2019], or transform synthetic images into photorealistic domain using adversarial training [Fried et al. 2019; Gecer et al. 2018]. Although these methods can modify facial region realistically, they are unable to complete the occluded parts, such as hair, body, and background.

More recently, a few methods propose to combine 3D prior with GAN to leverage its capacity in generating high-resolution photorealistic images of human faces [Deng et al. 2020; Geng et al. 2019; Xu et al. 2020]. The representative work StyleRig [Tewari et al. 2020a] provides a face rig-like 3D control over a pre-trained StyleGAN2 by mapping the control space of 3DMM to the latent space of StyleGAN2. Despite its success, it fails to preserve the visual quality when manipulating the real images. To address this issue, the following work PIE [Tewari et al. 2020b] proposes to embed the input image in the StyleGAN2 latent space ( $\mathcal{W}^+$  space) via hierarchical optimization. However, both methods use the entangled latent space of StyleGAN, thus cannot ensure the consistency of unchanged attributes and backgrounds (*e.g.* hair, clothes) after manipulation. Moreover, such methods cannot be applied to out-of-domain faces because building 3DMM for out-of-domain faces is either impractical or extremely expensive. In contrast, we propose Reduced StyleSpace, which has much better disentanglement among different face attributes and enables simultaneous control over them. We are also the first to generalize the rig-like control to out-of-domain face images via cross-domain adaption.

## 3 BACKGROUND AND OVERVIEW

### 3.1 Facial Attribute Parameterization

We adopt the widely-used 3D Morphable Model (3DMM) [Blanz and Vetter 1999] to parameterize the facial attributes, which provides compact semantic controls over both shape and texture.

In this work, we aim to support all the semantics that 3DMM can offer, *i.e.*, *geometry attributes* of identity  $\alpha$ , expression  $\beta$ , and pose  $T$ ; and *appearance attributes* of albedo  $\delta$  and illumination  $\gamma$ . Specifically,  $\alpha, \delta \in \mathbb{R}^{80}$  are the coefficients of morphable bases built through PCA,  $\beta \in \mathbb{R}^{64}$  are the coefficients corresponding to pre-defined expression bases,  $T$  is a rigid transformation containing both rotation  $r \in \mathbf{SO}(3)$  and translation  $t \in \mathbb{R}^3$ , and  $\gamma \in \mathbb{R}^{9 \times 3}$  are the order-3 Spherical Harmonics (SH) parameters for RGB channels. We denote this set of control parameters as  $P = (\alpha, \beta, T, \delta, \gamma)$  in the attribute space  $\mathcal{P}$ .

For the morphable model with  $N_v$  vertices, given the average geometry and appearance  $\{\bar{G}, \bar{A}\} \in \mathbb{R}^{3 \times N_v}$  and basis vectors  $\bar{G}^I, \bar{G}^E, \bar{A}$  for identity, expression, and albedo, respectively, we have the interpolated geometry  $G = \bar{G} + \alpha \times \bar{G}^I + \beta \times \bar{G}^E$  and appearance  $A = \bar{A} + \delta \times \bar{A}$ . Ultimately, the final positions  $V(P)$  and colors  $C(P)$  of all vertices are evaluated on each individual  $i$ -th vertex as follows:

$$V(P)_i = r \times G_i + t, \quad (1)$$

$$C(P)_i = A_i \cdot \sum_{b=1}^9 \gamma_b \cdot H_b(n_i), \quad (2)$$

where  $\{G_i, A_i\} \in \mathbb{R}^3$  are the vectors of the  $i$ -th vertex,  $H_b(n) \in \mathbb{R}$  is the response of the SH basis function  $H_b$  at normal direction of  $n$ , and  $\gamma_b \in \mathbb{R}^3$  is the illumination coefficients for the  $b$ -th SH band.

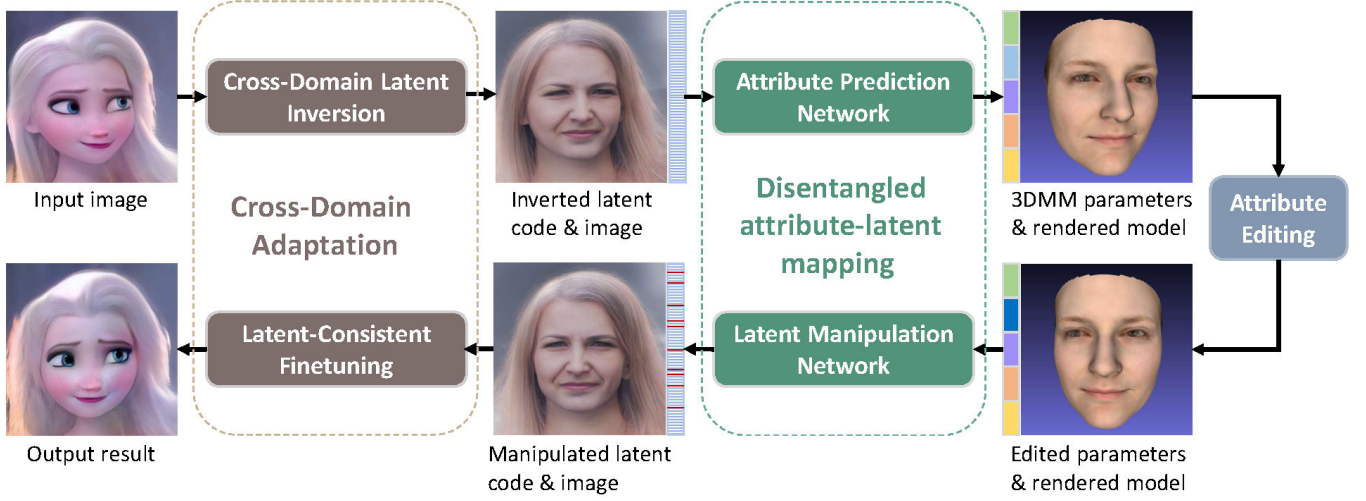


Fig. 2. The pipeline of the proposed method consists of two major stages: cross-domain adaptation and disentangled attribute-latent mapping. Given an input out-of-domain face image, the pipeline first reverts it into a latent code in the human face domain and then maps the latent code to 3DMM parameters via an attribute prediction network. Edits on pose, expression, illumination, etc., can be directly made to specific 3DMM parameters and then projected to disentangled latent code via a latent manipulation network. Finally, the manipulated latent code is fed into a StyleGAN2 finetuned with our latent-consistent finetuning techniques to generate the out-of-domain face, reflecting the corresponding edits.

### 3.2 Overview of Our Pipeline

Our cross-domain 3D-guided face manipulation pipeline is illustrated in Fig. 2. Given an input in-the-wild source image, the pipeline first applies domain-consistent latent inversion (§ 6.2) to optimize a latent code that can best reconstruct the image. Then, the pipeline updates the latent code (§ 5) to reflect various attribute manipulations controlled by 3DMM parameters (§ 3.1). This step basically consists of two parts: an attribute prediction network (§ 5.1) that predicts the 3DMM parameters from the latent code and a latent manipulation network (§ 5.2) that projects the parameter manipulation back onto the latent code. Both of them are built upon our reduced StyleSpace (§ 4.1) constructed through our attribute-adaptive layer selection method (§ 4.2) to ensure disentangled and precious controls. Finally, our space-consistent domain adaptation approach (§ 6.1) is adopted to map the manipulated latent code to an arbitrary out-of-domain face domain to achieve cross-domain face editing.

## 4 ATTRIBUTE-ADAPTIVE LATENT DECOMPOSITION

### 4.1 Reduced StyleSpace

A pre-trained StyleGAN2 [Karras et al. 2020] generator serves as a non-linear function  $G$  that maps the latent feature space  $\mathcal{W}^+$  to the image space:  $I = G(\mathbf{w})$ , where  $I$  is the generated image, and  $\mathbf{w} \in \mathbb{R}^{18 \times 512}$  is the latent vector. The raw  $\mathcal{W}^+$  space is a concatenation of 18 vectors of dimension 512, where each vector is transformed to channel-wise style parameters used to condition a corresponding AdaIN layer [Huang and Belongie 2017]. The space spanned by these style parameters is referred to as the *StyleSpace* [Wu et al. 2020], denoted as  $\mathcal{S}$ . For the target generation resolution of  $1024 \times 1024$ , a latent code in  $\mathcal{S}$  consists of style parameters for 26 layers with a total dimension of 9088, in which 17 layers of dimension 6048

apply to the feature maps, while the other 9 tRGB layers takes the remaining dimension of 3040.

Instead of  $\mathcal{W}^+$ , we adopt this *StyleSpace*  $\mathcal{S}$  as it is demonstrated to be more disentangled than other latent spaces including  $\mathcal{W}^+$  on embedding semantic properties [Wu et al. 2020], i.e., changing one attribute tends to be reflected on isolated variations in a few channels of  $\mathcal{S}$  without significant overlapping. Taking advantage of this desired property, we propose to further decompose the whole  $\mathcal{S}$  space into a certain number of subspaces, with each corresponding to one specific facial attribute. We call this set of subspaces the *Reduced StyleSpace*, where each attribute  $P^i$  is matched to a subspace  $\mathcal{S}^i \subseteq \mathcal{S}$ . Therefore, given a complete latent code  $\mathbf{w} \in \mathcal{S}$ , when manipulating one attribute  $P^i$  of its corresponding image, only part of the code  $\mathbf{w}^i \in \mathcal{S}^i$  should be altered accordingly. With such an attribute-adaptive reduced space, we have a compact and disentangled attribute embedding in the latent space, which not only facilitates the training performance but also makes it possible to edit one attribute without affecting irrelevant content or manipulate multiple attributes simultaneously without mutual interference.

### 4.2 Attribute-Adaptive Layer Selection

To faithfully select latent layers for each attribute, our intuition is to have an *agent* that maps the variation on each individual attribute to different layers in the space. We select those layers with strong responses since they tend to have strong correlations with the semantics controlled by that attribute.

The best choice for such a mapping agent should be driven by tasks that are relevant to encoding the 3DMM attributes to the latent space. Hence, we adopt the original StyleRig network trained on the *StyleSpace*  $\mathcal{S}$  as the agent, which demonstrates high-quality attribute-to-latent encoding. Here we denote the StyleRig network



**ALGORITHM 1:** Attribute-Adaptive Layer Selection Algorithm

---

**Input** : Pre-trained StyleRig network  $S$  on the  $\mathcal{S}$  space;  
 Data pairs  $\{(\mathbf{w}_1, P_1), (\mathbf{w}_2, P_2), \dots, (\mathbf{w}_N, P_N)\}$ ;  
**Output** : Selected layers  $\mathbb{L}^i$  for attribute  $\mathcal{P}^i$ ;  
**Define** :  $\text{REP}(P_a, P_b, i) = \{P_a^1, \dots, P_a^{i-1}, P_b^i, P_a^{i+1}, \dots, P_a^{|\mathcal{P}|}\}$ ;  
 $d \leftarrow 3$ ;  
 $\mathbb{L}^i \leftarrow \emptyset$ ;  
 $\Delta \mathbf{w} \leftarrow \sum_{a=1}^N \sum_{b=1}^N |S(\mathbf{w}_a, \text{REP}(P_a, P_b, i)) - \mathbf{w}_a| |P_a - P_b| / N^2$ ;  
**while**  $|\mathbb{L}^i| < 9$  **do**  
   **for**  $l \in [1, N_l]$  **do**  
    **for**  $c \in [1, N_c^l]$  **do**  
     **if**  $\Delta \mathbf{w}_{(l,c)} > d$  **then**  $n_{(l,c)} \leftarrow 1$ ;  
     **else**  $n_{(l,c)} \leftarrow 0$ ;  
   **end**  
   **if**  $\sum_{c=1}^{N_c^l} n_{(l,c)} > 0.25 * N_c^l$  **then**  $\mathbb{L}^i \leftarrow \mathbb{L}^i \cup \{l\}$ ;  
**end**  
 $d \leftarrow 0.95 * d$ ;  
**end**

---

simply as a function  $\mathbf{w}_g = S(\mathbf{w}_s, P_t)$ , which transforms the source code  $\mathbf{w}_s$  to  $\mathbf{w}_g$  controlled by the target attribute parameters  $P_t$ .

We propose the selection algorithm as detailed in Alg. 1, which eventually select a set of layers  $\mathbb{L}^i \in \mathcal{S}$  for each attribute  $\mathcal{P}^i$  to construct the corresponding subspace  $\mathcal{S}^i$ . Specifically, the correlation between a particular attribute and each layer is measured through the relative variability of latent value  $\mathbf{w}$  regarding changes of attribute  $P$ . Specifically, for the  $i^{\text{th}}$  parameter, given one random attribute-latent pair  $(\mathbf{w}, P)$ , we randomly change  $P^i$  and calculate the corresponding changes of  $\Delta \mathbf{w}$  using a pre-trained StyleRig network.  $\Delta \mathbf{w}$  is averaged on a large fixed set of random samples to ensure reliability. We select layers using a threshold for  $\Delta \mathbf{w}$ , and gradually decrease this threshold in case too few layers are selected. In our experiments, the same threshold is used for every parameter/layer, and is quite robust.

## 5 DISENTANGLED FACE MANIPULATION

At the core of our face manipulation framework is a mutual mapping structure between the *latent space*  $\mathcal{S}$  and the *attribute space*  $\mathcal{P}$ .

### 5.1 Attribute Prediction Network

Given a source latent code  $\mathbf{w}_s \in \mathcal{S}$  as the input to our framework, to manipulate the face encoded by this code, the first step is to interpret it to a set of 3D facial attributes  $P_s$  that are understandable and controllable. We achieve this through an *attribute prediction network*  $P = T(\mathbf{w})$ , which regresses 3DMM parameters  $P_s$  out of  $\mathbf{w}_s$  to best fit the subject face encoded in the source code. This translation network is also applied to the target latent code  $\mathbf{w}_t$  in the reference-based manipulation mode.

### 5.2 Latent Manipulation Network

After the source attribute parameters  $P_s$  are predicted with the attribute prediction network, user manipulation is performed in the attribute space, which either changes the parameters from the source  $P_s$  to the target  $P_t$  by incrementally editing one or some

of the attributes (edit-based manipulation) or entirely transferring from a reference (reference-based manipulation). The purpose of this attribute-to-latent space translation is then to update the latent code to inject the desired manipulation operations.

Thanks to the vector arithmetic properties [Radford et al. 2016] of the latent space embedding, we can potentially update the latent code via projecting the attribute changes to a linear displacement vector, of which the unit direction controls the identified editing semantics and the magnitude represents the manipulation intensity. This attribute-to-latent translation can be generally formulated as a mapping function  $\Delta \mathbf{w} = F(\mathbf{w}_s, P_s, P_t)$ , and the updated latent code to be  $\mathbf{w}_g = \mathbf{w}_s + \Delta \mathbf{w}$ .

The RigNet proposed in StyleRig [Tewari et al. 2020a] also falls into this formulation. Specifically, it contains an encoding structure  $E'$  that transforms the source code to feature vectors, and a decoding structure  $D'$  that takes both the feature vector and the target parameters to generate the final latent displacement vector:  $\Delta \mathbf{w} = F'(\mathbf{w}_s, P_t) = D'(E'(\mathbf{w}_s), P_t)$ , where  $\mathbf{w} \in \mathcal{W}^+$ . However, despite its huge success in achieving promising parametric 3D controls, there exist three types of *entanglements* that prevent more precise and flexible manipulation:

- The source latent code  $\mathbf{w}_s$  and the absolute target parameters  $P_t$  partially share the identity information of the subject;
- Correlations between different parameters in  $P_t$  prevent isolated manipulation through the translation networks;
- The  $\mathcal{W}^+$  space is not well disentangled regarding the semantics, editing one attribute could potentially affect the others.

In light of these issues, we adopt a *latent manipulation network* to achieve disentangled attribute-to-latent mapping, as shown in Figure 3. By using our reduced *StyleSpace* that decomposes  $\mathcal{S}$  into attribute-adaptive subspaces  $\{\mathcal{S}^1, \mathcal{S}^2, \dots, \mathcal{S}^{|\mathcal{P}|}\}$ , instead of jointly manipulating the entire target latent space  $\mathcal{W}^+$ , our framework consists of independent small translation networks  $\{F^1, F^2, \dots, F^{|\mathcal{P}|}\}$ , with each handles a specific target subspace of  $\mathcal{S}^i$  corresponding to an attribute  $P^i$  decided by our reduced *StyleSpace*. In addition, rather than manipulating the latent code on the absolute target parameters  $P_t$ , we feed the relative changes of parametric edits  $\Delta P = P_t - P_s$  to enforce the focus on the manipulation itself instead of the original identity. Our final attribute-to-latent translation framework is formulated as:

$$\Delta \mathbf{w}^i = F^i(\mathbf{w}_s^i, P_t^i - P_s^i) = D^i(E^i(\mathbf{w}_s^i), P_t^i - P_s^i), \quad (3)$$

$$\mathbf{w}_g^i = \mathbf{w}_s + \Delta \mathbf{w}^i. \quad (4)$$

Here, each network  $F^i$  consists of an encoder  $E^i$  and a decoder  $D^i$ .

### 5.3 Training Strategy

*Data Preparation.* The training data we use to train the system generally consists of tuples  $(\mathbf{w}, I, P)$ , in which  $\mathbf{w} \in \mathcal{S}$  is the latent code in the *StyleSpace*,  $I = G(\mathbf{w})$  is the image generated with StyleGAN2 generator  $G$  corresponding to the code  $\mathbf{w}$ , and  $P \in \mathcal{P}$  is the 3DMM parameters of the subject in  $I$ .

We first independently sample 5 latent vectors of size 512 in the  $\mathcal{W}$  space from a normal distribution. Then to construct the final latent code  $\mathbf{w}' \in \mathcal{W}^+$  of dimension  $18 \times 512$ , which is essentially a concatenation of 18 size-512 vectors, we randomly select one

from the 5 latent vectors for 18 times in a row and stack them together. After that, we convert  $\mathbf{w}'$  to  $\mathbf{w} \in \mathcal{S}$  and generate the face image  $\mathbf{I} = \mathbf{G}(\mathbf{w})$ . This technique is inspired by the mixing regularizer [Karras et al. 2020], which helps prevent the generator from assuming that the adjacent latent codes are correlated and improve the quality of the synthesized images. Finally, We predict the 3D parameters  $\mathbf{P}$  from the image  $\mathbf{I}$  using the off-the-shelf 3D face reconstruction method [Deng et al. 2019].

*Pre-Training for Attribute Prediction.* To avoid over-complicating the main training phase, we pre-train the attribute prediction network  $\mathbf{T}$  in a supervised manner with groundtruth pairs  $(\mathbf{w}, \mathbf{P})$ . In consideration of the varying contribution to the final 3D model of different parameters due to their distinct nature and strong correlation, we assign a weight for each individual parameter which is dynamically adjusted during the training. Specifically, we adopt the weighted parameter distance cost (WPDC) loss [Zhu et al. 2016]:

$$\mathcal{L}_{wpdc} = \|\epsilon \cdot (\mathbf{P} - \mathbf{T}(\mathbf{w}))\|_2^2, \quad (5)$$

where the parameter weight  $\epsilon$  is defined as the shape/color deviation caused solely by this parameter at the current value  $\mathbf{T}(\mathbf{w})$  if all other parameters are replaced with the groundtruth  $\mathbf{P}$ :

$$\begin{aligned} \epsilon_i &= \|\mathbf{V}(\text{REP}(\mathbf{P}, \mathbf{T}(\mathbf{w}), i)) - \mathbf{V}(\mathbf{P})\|_2 \\ &+ \|\mathbf{C}(\text{REP}(\mathbf{P}, \mathbf{T}(\mathbf{w}), i)) - \mathbf{C}(\mathbf{P})\|_2. \end{aligned} \quad (6)$$

Here REP is the replacement function that replaces the  $i$ -th element of  $\mathbf{P}_a$  with  $\mathbf{P}_b$ , as defined in Alg. 1.

*Self-Supervised Training for Latent Manipulation.* As ground-truth pairs are generally unavailable for in-the-wild images regarding most semantic manipulation operations, we adopt a self-supervised training solution to train the latent manipulation network.

Following StyleRig [Tewari et al. 2020a], training is performed on latent code pairs  $\{\mathbf{w}_s, \mathbf{w}_t\} \in \mathcal{S}$  together with their tuples  $(\mathbf{w}_s, \mathbf{I}_s, \mathbf{P}_s)$  and  $(\mathbf{w}_t, \mathbf{I}_t, \mathbf{P}_t)$ . The basic idea is to let the network manipulate the source code  $\mathbf{w}_s$  by replacing one single attribute of  $\mathbf{P}_s^i$  to  $\mathbf{P}_t^i$ . Through the networks  $\{\mathbf{T}, \mathbf{E}, \mathbf{D}\}$ , the generated latent code is calculated as  $\mathbf{w}_g = \mathbf{w}_s + \mathbf{D}(\mathbf{E}(\mathbf{w}_s), \text{REP}(\mathbf{T}(\mathbf{w}_s), \mathbf{T}(\mathbf{w}_t), i) - \mathbf{T}(\mathbf{w}_s))$ , with its parameters as  $\mathbf{P}_g = \mathbf{T}(\mathbf{w}_g)$ , where REP is the replacement function defined in Alg. 1. Given  $\mathbf{w}_g$ , we can enforce various self-supervision to ensure the corresponding image  $\mathbf{I}_g = \mathbf{G}(\mathbf{w}_g)$  conforms to the manipulation  $\mathbf{P}_s^i \rightarrow \mathbf{P}_t^i$ , while keeping other parameters  $\mathbf{P}_s^j, j \neq i$  unchanged. For the sake of simplicity, here we define two terms that are important to the self-supervisions:  $\mathbf{P}_{tg} = \text{REP}(\mathbf{P}_t, \mathbf{P}_g, i)$  and  $\mathbf{P}_{gs} = \text{REP}(\mathbf{P}_g, \mathbf{P}_s, i)$ , which should be as close as possible to the target and source parameters, respectively, to measure the accuracy and the quality of disentanglement of the manipulation.

First of all, we utilize the differentiable renderer to measure photometric similarities in the image space and back-propagate them to the attribute space. This rendering loss is defined as:

$$\mathcal{L}_{render} = \|\mathbf{R}(\mathbf{P}_{tg}) - \mathbf{M}(\mathbf{P}_t) \cdot \mathbf{I}_t\|_2^2 + \|\mathbf{R}(\mathbf{P}_{gs}) - \mathbf{M}(\mathbf{P}_s) \cdot \mathbf{I}_s\|_2^2, \quad (7)$$

where  $\mathbf{R}(\mathbf{P})$  is the differentially rendered face image from parameters  $\mathbf{P}$ , with its corresponding occupancy mask  $\mathbf{M}(\mathbf{P})$ , i.e., 2D regions covered by the projected morphable model. The dimensions of these images are made equal to  $\mathbf{I}_s$  and  $\mathbf{I}_t$ .

By using sparse landmarks that are pre-labeled on the mesh, we also adopt the landmark loss:

$$\begin{aligned} \mathcal{L}_{land} &= \|\Pi(\mathbf{L}(\mathbf{V}(\mathbf{P}_{tg}))) - \Pi(\mathbf{L}(\mathbf{V}(\mathbf{P}_t)))\|_2^2 \\ &+ \|\Pi(\mathbf{L}(\mathbf{V}(\mathbf{P}_{gs}))) - \Pi(\mathbf{L}(\mathbf{V}(\mathbf{P}_s)))\|_2^2, \end{aligned} \quad (8)$$

where the function  $\mathbf{L}(\mathbf{G}) \in \mathbb{R}^{N_l \times 3}$  samples the 3D landmark positions from mesh vertices  $\mathbf{G}$  and the function  $\Pi(\mathbf{L}) \in \mathbb{R}^{N_l \times 2}$  projects landmarks  $\mathbf{L}$  to 2D positions on the image plane with the weak-perspective camera model. We use  $N_l = 68$  in our experiments.

Similarly, we propose the contour loss to further improve the consistency of face shape:

$$\begin{aligned} \mathcal{L}_{cont} &= \|\mathbf{L}'(\mathbf{V}(\mathbf{P}_{tg})) - \mathbf{L}'(\mathbf{V}(\mathbf{P}_t))\|_2^2 \\ &+ \|\mathbf{L}'(\mathbf{V}(\mathbf{P}_{gs})) - \mathbf{L}'(\mathbf{V}(\mathbf{P}_s))\|_2^2, \end{aligned} \quad (9)$$

where  $\mathbf{L}'(\mathbf{G}) \in \mathbb{R}^{N_c \times 3}$  samples only a subset of landmarks that are around the face outer contour ( $N_c=17$ ). Different from the 2D landmark loss  $\mathcal{L}_{land}$  (Eq. 8), this contour loss is defined on 3D positions without transformation-dependent projection, which avoids the pose bias due to inaccurate rotation.

$$\mathcal{L} = \lambda_r \mathcal{L}_{render} + \lambda_l \mathcal{L}_{land} + \lambda_c \mathcal{L}_{cont}. \quad (10)$$

## 6 CROSS-DOMAIN ADAPTATION

Through our in-domain latent-attribute mapping pipeline, we are able to perform disentangled manipulation of real face images. In this section, we introduce how we adapt this 3D-guided face manipulation framework to handle cross-domain faces, including non-photorealistic paintings, cartoon portraits, animals, etc.

### 6.1 Latent-Consistent Finetuning

As introduced in § 5, based on the StyleGAN2 generator  $\mathbf{G}$  that is pre-trained on the real-face domain  $\mathcal{A}$ , we edit latent code in the latent space  $\mathcal{S}$  via its mutual mapping with the facial attribute space  $\mathcal{P}$  parameterized by human 3DMM. To make the best of the well-defined semantics parameterization on human faces and to maintain the consistency of our disentangled latent embedding, we propose to reuse both spaces for cross-domain face manipulation and perform domain-consistent generator finetuning to ensure that the cross-domain latent code could be effectively mapped to each domain-specific generator.

Specifically, for a new face domain  $\mathcal{B}$ , we train a generator  $\mathbf{G}^*$  on  $\mathcal{B}$  which shares a same latent embedding as  $\mathbf{G}$ . Given an source image  $\mathbf{I}_s \in \mathcal{B}$ , we apply cross-domain latent inversion to get latent code  $\mathbf{w}_s \in \mathcal{S}$  so that  $\mathbf{G}^*(\mathbf{w}_s) \sim \mathbf{I}_s$ . With  $\mathbf{w}_s$ , in-domain face manipulation is used to update the latent code to  $\mathbf{w}_g$  and produce the final image  $\mathbf{I}_g = \mathbf{G}^*(\mathbf{w}_g)$  in domain  $\mathcal{B}$ .

We sample training images from the target domain  $\mathcal{B}$  and use them to finetune the original real-face StyleGAN2 generator  $\mathbf{G}$  to get the domain-specific generator  $\mathbf{G}^*$  for domain  $\mathcal{B}$ . However, straightforward finetuning which updates all network parameters would distort the underlying latent space. To keep the latent consistency across all generators, we freeze the network layers relevant to latent embedding. To be specific, since we use the *StyleSpace*  $\mathcal{S}$ , all

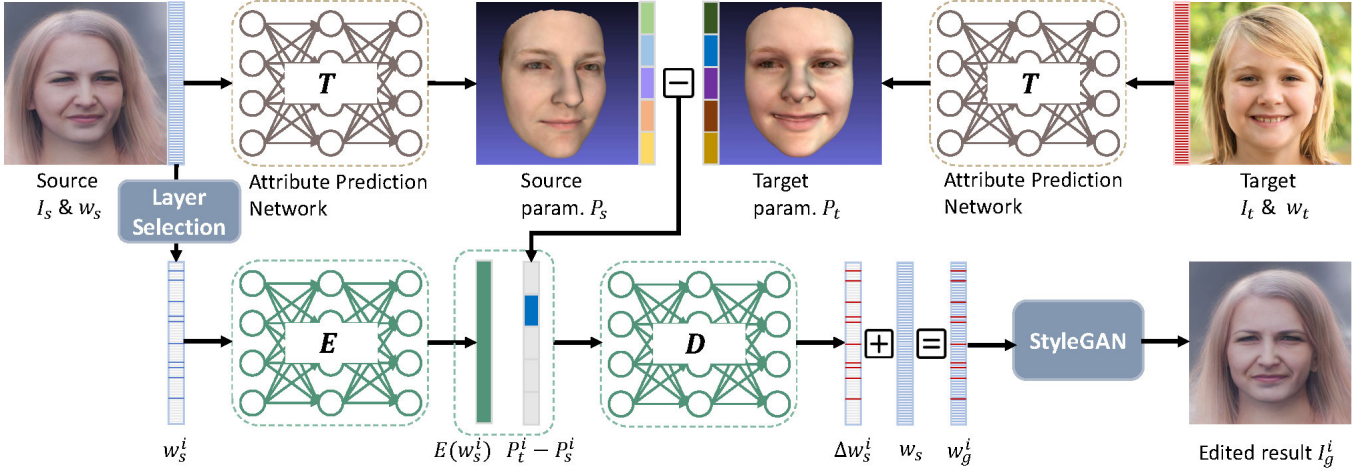


Fig. 3. Framework of disentangled face manipulation. Given a source latent code  $w_s$ , to manipulate the face  $I_s$  encoded by this code, the first step is to interpret it to a set of 3D facial attributes  $P_s$  through an attribute prediction network  $T$ . And then we take the difference of  $i$ -th attribute between  $P_s$  and the target parameter  $P_t$ , getting  $\Delta P^i = P_t^i - P_s^i$ .  $\Delta P^i$  together with the source code  $w_s^i$  in subspace  $S^i$  is fed into a latent manipulation network (composed of an encoder  $E$  and a decoder  $D$ ), to predict  $\Delta w^i$  for updating source latent code  $w_s$  and generating the edited result  $I_g^i$  with a pretrained StyleGAN.

style convolution layers and tRGB layers are fixed during the fine-tuning. Besides that, we use exactly the same loss functions and hyper-parameters as in StyleGAN2 [Karras et al. 2020]. Eventually, we have  $G$  and  $G^*$  that share a same latent space  $\mathcal{S}$ . A manipulated latent code  $w_g$ , conditioned on human-face 3DMM attributes, could be mapped to different domains with the same controlled semantics.

## 6.2 Cross-Domain Latent Inversion

To complete the pipeline that manipulates an in-the-wild image from an arbitrary face domain, the last piece of the puzzle we need is a cross-domain latent inversion component  $I$  that can robustly embed an in-the-wild image from the face domain  $\mathcal{B}$  into our latent space  $\mathcal{S}$ :  $w_s = I(I_s), w_s \in \mathcal{S}, I_s \in \mathcal{B}$ .

We take an optimization-based approach with perceptual and MSE loss between the original and the reconstructed image:

$$w'_s = \arg \min_{w'} (\|I_s - G^*(w')\|_2^2 + \sum_{l=1}^L \|\Psi_l(I_s) - \Psi_l(G^*(w'))\|_1), \quad (11)$$

where  $\Psi_l(I)$  computes the activation feature map of image  $I$  at the  $l$ -th selected layer of the VGG-19 network pre-trained on ImageNet. Here  $w' \in \mathcal{W}^+$ . We transform it to the  $\mathcal{S}$  space using the style convolutional layers  $S$  as  $w_s = S(w'_s)$ .

## 7 EXPERIMENTS

### 7.1 Implementation Details

**Training Parameters.** We train both the attribute prediction and the latent manipulation networks using Adam optimizer for 60 epochs with the batch size of 128. The initial learning rate is set to 0.01 and decayed by 0.1 every 10 epochs. For the latent manipulation network, the weights of each loss term are set as  $\lambda_r = 1$ ,  $\lambda_l = 0.001$ , and  $\lambda_c = 0.01$  for all our experiments, respectively. To train the domain-specific StyleGAN2 generator  $G^*$  for each out-of-domain domain, we finetune the pretrained StyleGAN2 for 32,000 iterations with the batch size of 16 on each dataset using the same learning rate

scheduler but a smaller learning rate of 0.002. Besides, Pytorch3D is used as the differential renderer.

**Network Architecture.** The attribute prediction network consists of 5 MLP layers with 4 ELU activations after each intermediate layer, and the hidden dimensions are 4096, 2048, 1024, 512, 257 respectively. For the latent manipulation network, the encoder  $E^i$  for each attribute  $i$  consists of  $N$  MLP layers, where  $N$  is set as 9, 4, 7, 9, 4 for expression, pose, albedo, illumination and shape respectively. And the decoder also consists of 3 MLP layers. Uniform weight initialization is used for the networks by default.

**Out-of-Domain Datasets.** We use five out-of-domain face datasets to demonstrate the cross-domain face manipulation effects, including Ukiyo-e Face [Pinkney 2020], AFHQ Dog [Choi et al. 2020], WikiArt Dataset<sup>1</sup>, Danbooru2018 [Anonymous et al. 2021], and Disney Face involving 400 online images of Disney cartoon characters collected by ourselves.

### 7.2 Ablation Study for Disentanglement

Disentangling different face attributes in the editing process is crucial for enabling the fine-grained controllability of each attribute and guaranteeing the final editing quality. To improve the disentanglement, our method has two key designs: using the reduced *StyleSpace* where one specific facial attribute has one specific subspace rather than the  $\mathcal{W}^+$  space used in [Tewari et al. 2020a], and using the relative changes of the parametric edits  $\Delta P$  rather than the absolute target parameters  $P_t$  in the latent manipulation network.

To quantitatively evaluate the disentanglement for each facial attribute, we randomly sample a set of source latent codes  $\{w_s^i\}_{i=1}^n$  from the latent space  $\mathcal{S}$  and measure the average 3DMM parameter difference before and after editing along the dimension of each attribute respectively. More specifically, for each latent code  $w_s^i$ , we

<sup>1</sup><https://github.com/cs-chan/ArtGAN/tree/master/WikiArt%20Dataset>

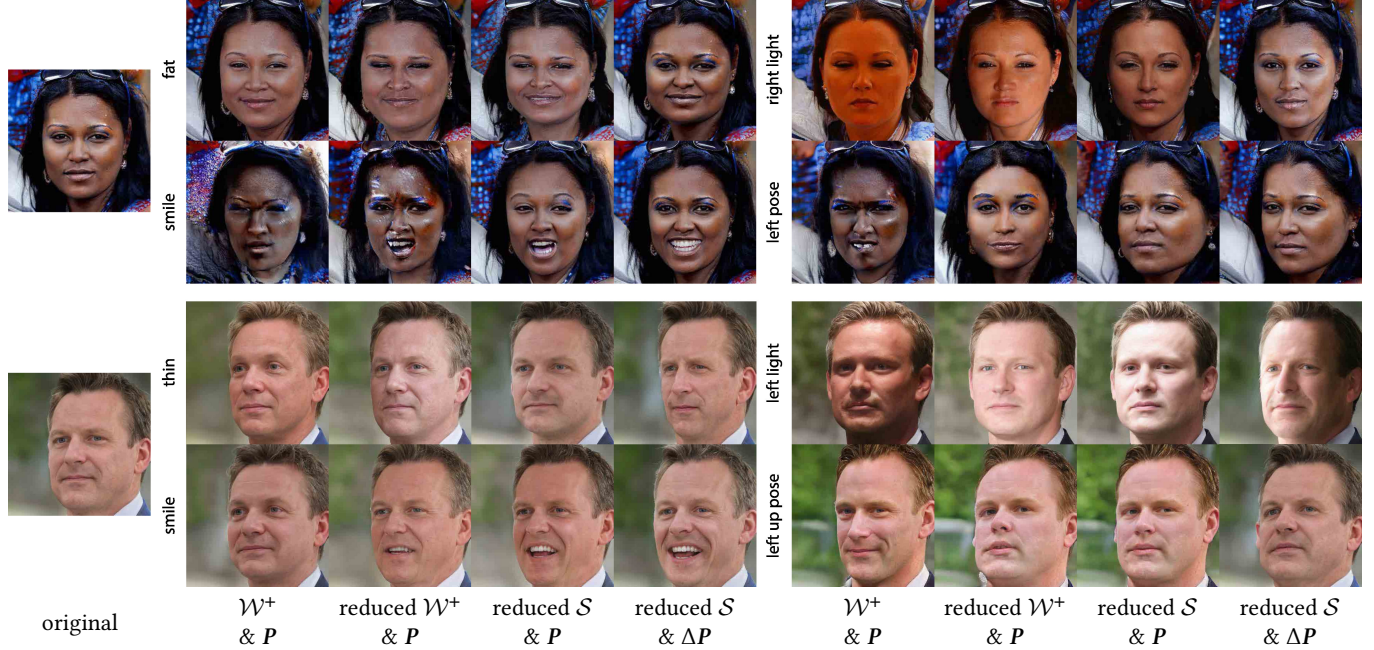


Fig. 4. Visual results of different spaces for the attribute prediction network and different inputs for the latent manipulation network. With our reduced *StyleSpace* and the relative parametric changes  $\Delta P$ , our results (last column of each case) achieve the best attribute disentanglement quality.

Table 1. Quantitative evaluation of disentanglement quality of attribute manipulation using different spaces ( $\mathcal{W}^+$  Space, reduced  $\mathcal{W}^+$  Space, *StyleSpace*, and reduced *StyleSpace*) for the attribute prediction network and different inputs ( $P$ ,  $\Delta P$ ) for the latent manipulation network. Lower is better.

Method	Shape	Exp	Illumination	Pose	Albedo	Avg.
$\mathcal{W}^+ \& P$	24.43	24.71	25.03	21.06	25.12	24.07
$S \& P$	23.91	24.13	24.87	16.79	24.97	22.93
$S \& \Delta P$	22.18	20.13	24.72	12.4	24.67	20.82
reduced $\mathcal{W}^+ \& P$	24.11	24.69	24.99	19.09	24.79	23.51
reduced $S \& P$	23.50	22.50	24.74	14.98	24.77	22.10
reduced $S \& \Delta P$	<b>19.62</b>	<b>18.75</b>	<b>24.70</b>	<b>8.81</b>	<b>24.22</b>	<b>19.22</b>

first map it into the corresponding 3DMM parameter  $P_s^i = T(w_s^i)$  via the attribute prediction network, then we do the editing in the 3DMM space and get the target edited 3DMM parameter  $P_t^i$ . Next,  $P_t^i$  is mapped back to the latent space via the latent manipulation network and get the edited latent code  $w_g^i$ . Finally,  $w_g^i$  is mapped to the 3DMM space to get the 3DMM parameter  $P_g^i$  for the actually edited face. The average L1 distance between  $P_g^i$  and  $P_t^i$  along the attribute-specific dimension will be regarded as the disentanglement metric for each attribute. Note values of different attributes cannot be compared directly, because value range and data variation of different attributes are different. That’s also why pose is captured significantly lower than all other attributes.

*Benefits of the Reduced StyleSpace.* To demonstrate the advantage of the newly proposed reduced *StyleSpace* over the  $\mathcal{W}^+$  space used in [Tewari et al. 2020a], we keep using  $P$  (not  $\Delta P$ ) and replace the reduced *StyleSpace* in our framework with the original  $\mathcal{W}^+$  space, where different attributes are not disentangled, and the variant “reduced  $\mathcal{W}^+$  space” where each attribute is disentangled and has its own subspace described in Section 4.1. As shown in Table 1, using independent subspaces for each attribute can improve the disentanglement metric from 24.07 to 23.51. By using the *StyleSpace* rather than  $\mathcal{W}^+$  Space, it can further improve the disentanglement metric into 22.10. We think the underlying reason of improvements may come from two aspects: 1) Using one independent subspace for each attribute can not only avoid the interference from other attributes but also reduce the learning difficulty from too high space dimension; 2) *StyleSpace* has better disentanglement than  $\mathcal{W}^+$  space inherently.

*Benefits of  $\Delta P$ .* Similarly, as shown in the last row of Table 1 (our default setting), by replacing the absolute target parameter  $P_t$  with the relative change  $\Delta P$  as the input of the latent manipulation network, the disentanglement metric significantly improves from 22.10 to 19.22. And a similar observation by comparing  $S \& P$  and  $S \& \Delta P$ . One possible explanation is that, since  $P_t$  contains other attributes’ information, forcing the manipulation network to deliberately ignore such information (if given) will increase the learning difficulty of one specific attribute. In contrast, using  $\Delta P$  is equivalent to explicitly telling the network which attribute should be changed and which attribute should not be changed. Intuitively, it is also more natural as the design motive of the manipulation network is just to



build the relationship between the change in the latent space and the change in the 3DMM space.

Besides the above quantitative results, we also provide some representative visual results in Figure 4. Compared to our default setting, the three remaining settings either produce serious artifacts or incur unexpected attribute change (e.g., identity) when editing one specific attribute. Taking the “right light” editing case of the woman image as an example, the “identity” attribute seems changed in all the baseline settings. This is consistent with the conclusion drawn from the quantitative evaluations.

### 7.3 Comparisons

To demonstrate our advantages in cross-domain face manipulation, we conduct comparisons with different types of representative baselines, including 2D warping based methods: Moving-Least-Squares (MLS) deformation [Schaefer et al. 2006] and first order motion model [Siarohin et al. 2019b], latent space manipulation based methods: GANSpace [Härkönen et al. 2020] and InterFaceGAN [Shen et al. 2020], and 3DMM guided methods: StyleRig [Tewari et al. 2020a] and PIE [Tewari et al. 2020b].

*Comparisons with 2D Warping Based Methods.* We first compare our method with 2D warping based methods. In Figure 5, we show the editing results in terms of shape, pose, and expression changes, respectively. For MLS, we use target landmarks to guide the deformation. For first-order motion, we construct video sequences from target images to drive the non-human face. Since pure pixel-level warping method like MLS cannot synthesis originally occluded regions, it fails at manipulations such as mouth opening. While the first-order motion model can handle occlusion quite well, it suffers from the same issue as MLS that 2D warping field interpolated from sparse landmarks cannot well represent 3D deformation such as head rotation. In contrast, by leveraging the generative property of GAN and the 3D prior guidance, our method achieves higher fidelity and controllability without these issues.

*Comparisons with Latent Space Manipulation Methods.* We then compare our method with GANSpace and InterFaceGAN that perform semantic image control via latent space manipulation in Figure 6. For these manipulation methods, the pretrained StyleGAN2 model is finetuned on the same non-human face dataset like ours. Since GANSpace is an unsupervised method without explicit supervision of the semantic attributes, it is difficult to locate the corresponding latent variable for a given attribute. Also, it is not guaranteed that all the factors of interest will be disentangled, e.g., the hair is also changed when “opening mouth”. For InterFaceGAN, as their original version only supports human faces, we combine it with our cross-domain adaptation method when applied to out-of-domain faces. Compared to GANSpace, InterFaceGAN allows explicit controls with supervision, but the required binary labeling hardly exists for some attributes, thus not supporting the controllable editing for such attributes (e.g., changing the illumination or shape). In contrast, our method can disentangle different attributes much better and produce better editing quality.

*Comparisons with 3DMM Guided Methods.* Finally, we compare our method with StyleRig and PIE, which also leverage the 3DMM

Table 2. User study results. The values for each attribute represent the user preference rate (%), the higher, the better) by comparing the editing results among different methods. “-” means that the method does not support editing this attribute or the editing direction is not available in their paper.

Method	Shape	Exp	Illumination	Pose	Avg.
GANSpace	10.57	0.86	13.71	0.57	6.43
First Order	-	4.29	-	2.86	3.58
InterFaceGAN	-	2.86	-	2.00	2.43
MLS	0.57	0.57	-	-	0.57
Ours	<b>88.86</b>	<b>91.42</b>	<b>86.29</b>	<b>94.57</b>	<b>90.29</b>

guidance for face manipulation. But due to the lack of 3DMM parametric space for out-of-domain faces, they do not support out-of-domain faces editing inherently. Therefore, we only do the comparison on real face images in Figure 7. For a fair comparison, we directly fetch the results of StyleRig and PIE from the authors’ project page. Without our reduced latent space and disentangled attribute-latent mapping, on the one hand, they cannot manipulate identity-related attributes such as shape and albedo. On the other hand, they often fail at editing a single attribute without affecting other content or manipulating multiple attributes simultaneously. For example, when changing the expression in the first case, the glass structure is a little twisted. Similarly, when changing the expression in the second case, the albedo and background are also significantly changed with obvious artifacts in the edited result.

### 7.4 User Study

To quantitatively compare our method and baseline methods in terms of the out-of-domain face editing quality, we conduct a user study and let user choose their preferred results. Since StyleRig and PIE only support real face image editing, we do not include them here. Similarly, as changing the albedo attribute is not supported in the remaining baseline methods, only the controllable editing results for shape, expression, illumination, and pose are compared. If the baseline method does not support editing one specific attribute, it will be ignored in that corresponding comparison. Specifically, for each controllable attribute editing comparison of one domain, we randomly choose a total of 10 images. Thus for 5 domains with 4 attributes, we have 200 images in total. Then we use different methods to edit the target attribute while keeping other attributes unchanged. Then for each case, we will show the edited results together with the original image to total 35 participants and ask them the question “Which result is the best for the “X change” while keeping other attributes in the original image intact, including the identity, hair style and other face details?” Here, “X change” indicates different attribute editing operation, like “being fat” and “turn right”. The final preference rate is defined as the averaged percentage of one specific method is selected as the best.

As shown in Table 2, the users prefer the editing results generated from our method than all the baseline methods by a large margin, which is consistent to the conclusion drawn by the above visual comparisons.

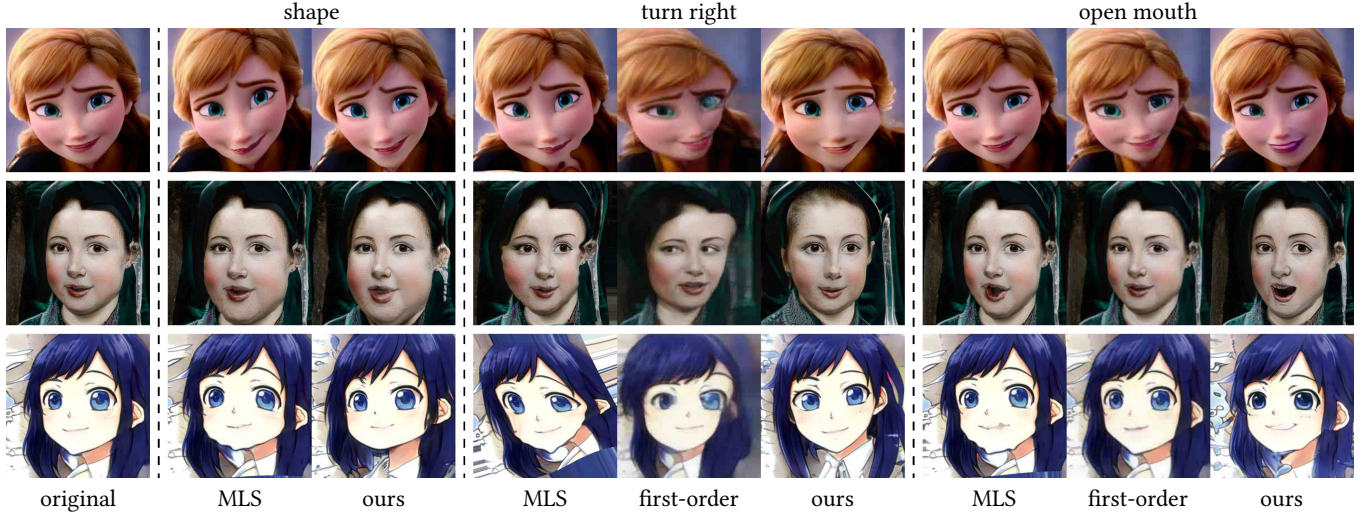


Fig. 5. Comparison results with 2D warping-based methods MLS and First Order. Since First Order does not support shape change, we do not include its result for the left “Be Thin/Fat” case. It can be seen that our method can achieve much better editing results than both MLS and First Order, which either cause obvious artifacts or fail to achieve the target editing effect (e.g., “Open mouth” for First Order).

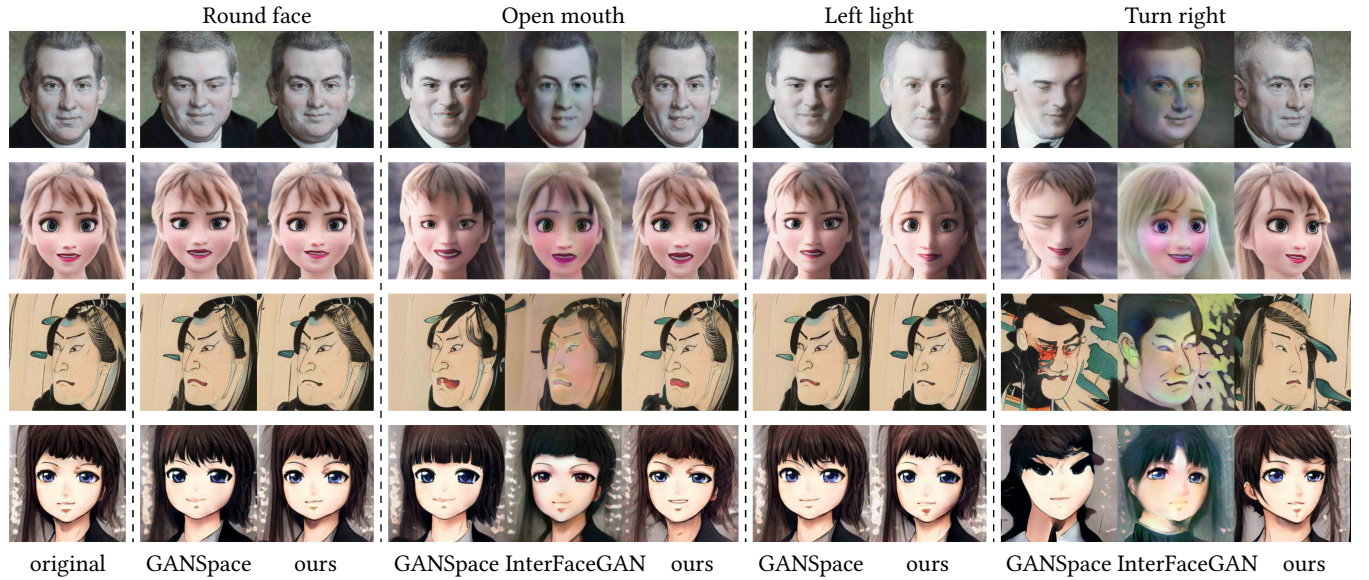


Fig. 6. Comparison results with GAN latent space manipulation methods: GANSpace and InterFace. Since InterFaceGAN does not support shape change and illumination change, we do not provide its results for the first and third case.

## 8 RESULTS

### 8.1 User Manipulation InterFace

To facilitate the editing process for users, we develop an interactive editing system shown in Figure 8. To edit one out-of-domain face, 1) the system first inverts this image into a latent code by using the finetuned StyleGAN2; 2) then map this latent code into the 3D parameters by using the attribute prediction network and show the 3D mesh in the main window; 3) the users can edit images

by adjusting values for different 3DMM bases of each attribute in the left panel or directly edit the 3D geometry by dragging the 3D facial landmarks, and the backend will use the Laplacian mesh deformation algorithm to get a well-edited 3D mesh; 4) the edited 3D parameters will be mapped back to the latent space by using the latent manipulation network; 5) finally feed the edited latent code into the finetuned StyleGAN2 to get the final editing result.



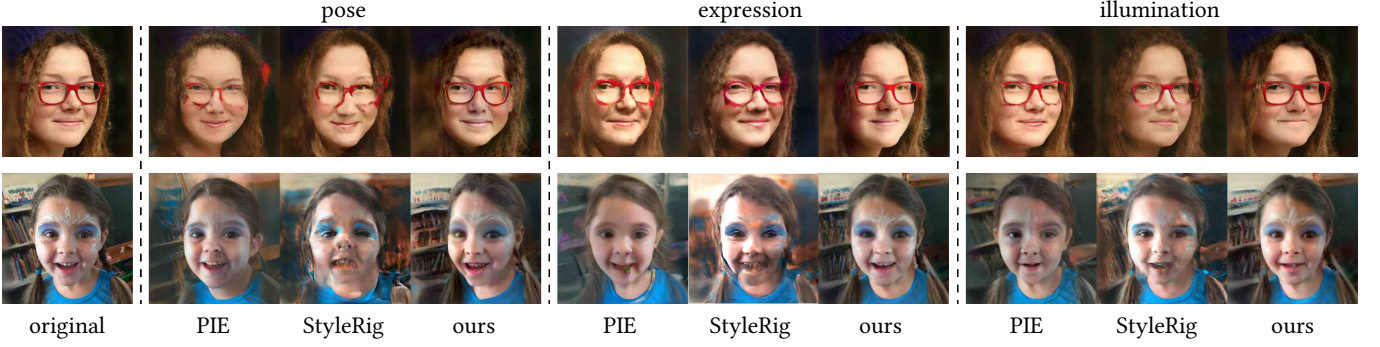


Fig. 7. Comparison with 3DMM guided methods StyleRig and PIE. As they do not support out-of-domain face editing, we only show the editing results for the real faces. It shows that our method can disentangle different attributes better and produce more controllable editing results with higher quality.

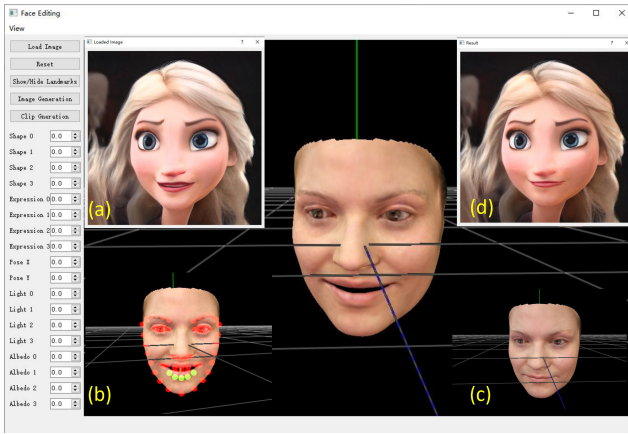


Fig. 8. Our interactive cross-domain face editing system. (a) shows the input out-of-domain face image. (b) shows the windows that users can edit the 3D mesh of the corresponding real face in a fine-grained way by controlling the facial landmarks. (c) shows the 3D mesh editing result. (d) is the final generated result. The left panel shows the editable parameters for each facial attribute, including shape, expression, pose, illumination and albedo.

## 8.2 Multi-Attribute Manipulation

In the above comparison, all the results are obtained by editing one single attribute while keeping other attributes unchanged. But in fact, it is also very easy for our method to support editing multiple attributes simultaneously. Specifically, we directly add all the  $\Delta w^i$  onto the original latent code  $w_s$  for each attribute  $i$  involved and then feed the edited latent code into the finetuned StyleGANv2 to generate the final result. In Figure 9, we provide three representative results of editing two to four attributes simultaneously. But for the baseline method StyleRig and PIE, they only support editing multiple attributes in an incremental way (one by one). The underlying reason is that they cannot disentangle each attribute very well and the latent code change for different attributes will interfere with each other. In contrast, our disentangled attribute-latent mapping greatly reduces the overlapping between attributes, which

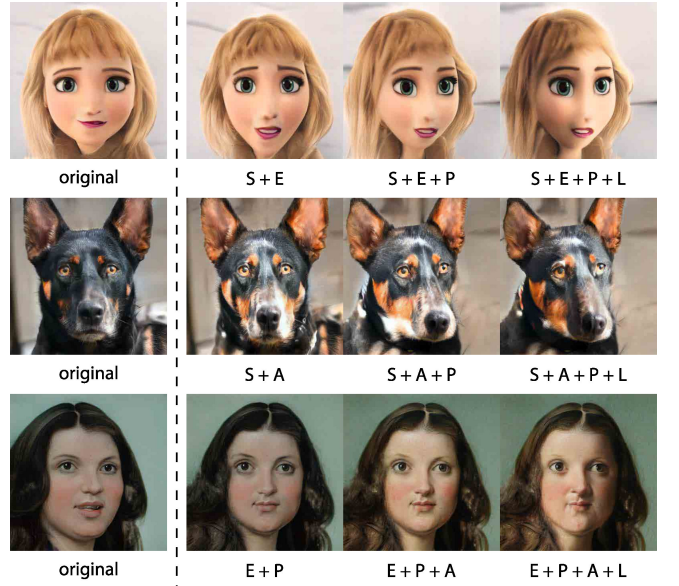


Fig. 9. Multi-attribute manipulation results for shape(S), expression(E), pose(P), illumination(I) and albedo(A). Our method can manipulate multiple attributes simultaneously and generate visually plausible results.

makes it possible to achieve multi-attribute manipulation without introducing significant interference.

## 8.3 Continuous Manipulation

Besides the above multi-attribute manipulation, our method also naturally supports continuous manipulation for changing the editing strengths of one specific attribute and changing one specific attribute to another attribute in a smooth way. Figure 11 has shown three such continuous editing results. Taking the third baby girl as an example, we first gradually change the light from the right to the middle, and then switch to expression change and change the mouth status from close to open. This continuous editing function not only demonstrates the great disentanglement of our method

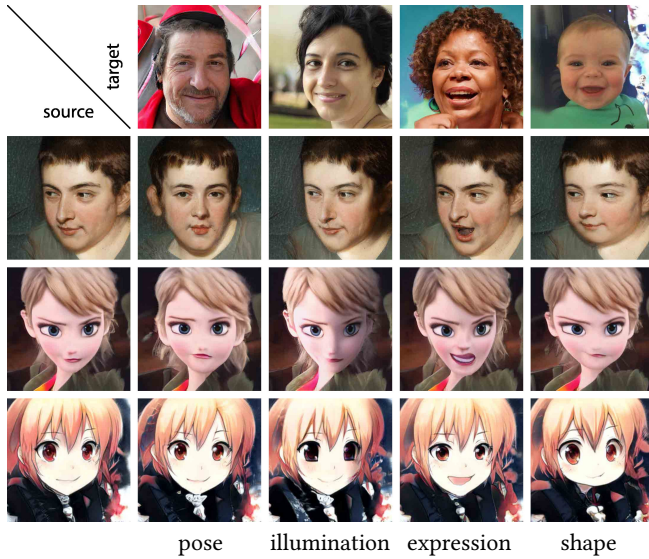


Fig. 10. Visual results to show the cross-domain generalization ability. Given the same reference image, we can edit different out-of-domain domain images in terms of one specific attribute while keeping other attributes unchanged.

but also potentially enables one new effect, i.e., bringing one still out-of-domain face image to life.

#### 8.4 Cross-Domain Generalization

In Figure 10, we further demonstrate the cross-domain generalization ability of our method. In details, for the source images from different out-of-domain face domains, we aim to change the same attribute of all of these images by adopting the same reference face image. It shows that, for each attribute specified by the real face, our method can consistently follow the editing direction and achieve the corresponding editing effect. This indeed echoes our main contribution, i.e., disentangled out-of-domain face manipulation via bridging the discrepancy between real human face and out-of-domain face in the semantic latent space.

## 9 CONCLUSIONS

We have presented the first approach for semantic attribute manipulation of out-of-domain faces via adopting 3DMM of human faces as the proxy. To this end, we devised a cross-domain adaptation method that bridges domain discrepancies and allows 3DMM parameter edits on the human face to be faithfully reflected on the out-of-domain face image. In addition, we proposed a reduced latent space and disentangled attribute-latent mapping to guarantee disentangled and precise controls for each semantic attribute. With our approach, there is no need to build 3DMM for a specific out-of-domain face domain, and intuitive parameter editing, including the head pose, facial shape, expression, albedo, and illumination, is well supported for arbitrary out-of-domain faces. The visually pleasant quality and user-friendly control show the great potential

of our method for many exciting applications in the areas of design, cartoons, animations, and games.

Since our approach is the first step towards 3DMM-based manipulation of out-of-domain faces, there is still room for further improvement. First, some fine-grained edits beyond the expressivity of 3DMM, like wrinkling nose or wearing glasses, are not supported by our method. This might be improved by adopting a high-fidelity 3DMM. Second, some 3DMM attribute edits may not be effective in some specific domains. As shown in Figure 12, we cannot control a dog to open mouth or change the illumination of a Ukiyo-e face since their datasets for fine-tuning the StyleGAN2 do not contain such variations. Third, if our method is directly applied to video frames, some popping artifacts will be noticeable. This can be alleviated by adding temporal coherency loss [Chen et al. 2017] when training our disentangled attribute-latent mapping network. Moreover, our method needs to fine-tune the StyleGAN2 in a specific domain before manipulating images of that domain. In order to increase the scalability of our method, how to fast update a StyleGAN2 with small-scale data or even a single shot would be a worth-exploring direction in the future.

## REFERENCES

- Rameen Abdal, Yipeng Qin, and Peter Wonka. 2019. Image2StyleGAN: How to Embed Images Into the StyleGAN Latent Space?. In *ICCV 2019*. 4431–4440.
- Rameen Abdal, Peihao Zhu, Niloy J. Mitra, and Peter Wonka. 2020. StyleFlow: Attribute-conditioned Exploration of StyleGAN-Generated Images using Conditional Continuous Normalizing Flows. *CoRR* abs/2008.02401 (2020).
- Anonymous, Danbooru community, and Gwern Branwen. 2021. Danbooru2020: A Large-Scale Crowdsourced and Tagged Anime Illustration Dataset. <https://www.gwern.net/Danbooru2020>. <https://www.gwern.net/Danbooru2020> Accessed: DATE.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein Generative Adversarial Networks. In *ICML 2017*, Vol. 70. 214–223.
- Hadar Averbuch-Elor, Daniel Cohen-Or, Johannes Kopf, and Michael F. Cohen. 2017. Bringing portraits to life. *ACM Trans. Graph.* 36, 6 (2017), 196:1–196:13.
- Volker Blanz and Thomas Vetter. 1999. A Morphable Model for the Synthesis of 3D Faces. In *SIGGRAPH 1999*. 187–194.
- Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. 2017. Coherent Online Video Style Transfer. In *ICCV 2017*. 1114–1123.
- Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. 2020. StarGAN v2: Diverse Image Synthesis for Multiple Domains. In *CVPR 2020*. 8185–8194.
- Kevin Dale, Kalyan Sunkavalli, Micah K. Johnson, Daniel Vlasic, Wojciech Matusik, and Hanspeter Pfister. 2011. Video face replacement. *ACM Trans. Graph.* 30, 6 (2011), 130.
- Yu Deng, Jialong Yang, Dong Chen, Fang Wen, and Xin Tong. 2020. Disentangled and Controllable Face Image Generation via 3D Imitative-Contrastive Learning. In *CVPR 2020*. 5153–5162.
- Yu Deng, Jialong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. 2019. Accurate 3D Face Reconstruction With Weakly-Supervised Learning: From Single Image to Image Set. In *CVPR Workshops 2019*. 285–295.
- Ohad Fried, Eli Shechtman, Dan B. Goldman, and Adam Finkelstein. 2016. Perspective-aware manipulation of portrait photos. *ACM Trans. Graph.* 35, 4 (2016), 128:1–128:10.
- Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B. Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. 2019. Text-based editing of talking-head video. *ACM Trans. Graph.* 38, 4 (2019), 68:1–68:14.
- Pablo Garrido, Levi Valgaerts, H. Sarmadi, I. Steiner, Kiran Varanasi, Patrick Pérez, and Christian Theobalt. 2015. VDub: Modifying Face Video of Actors for Plausible Visual Alignment to a Dubbed Audio Track. *Comput. Graph. Forum* 34, 2 (2015), 193–204.
- Baris Gecer, Binod Bhattarai, Josef Kittler, and Tae-Kyun Kim. 2018. Semi-supervised Adversarial Learning to Generate Photorealistic Face Images of New Identities from 3D Morphable Model. In *ECCV 2018*, Vol. 11215. 230–248.
- Jiahao Geng, Tianjia Shao, Youyi Zheng, Yanlin Weng, and Kun Zhou. 2018. Warp-guided GANs for single-photo facial animation. *ACM Trans. Graph.* 37, 6 (2018), 231:1–231:12.
- Zhenglin Geng, Chen Cao, and Sergey Tulyakov. 2019. 3D Guided Fine-Grained Face Manipulation. In *CVPR 2019*. 9821–9830.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. 2014. Generative Adversarial





Fig. 11. Examples of continuous manipulation on the attributes, including shape(S), expression(E), pose(P), illumination (I), and albedo(A). And PX and PY are the horizontal and vertical pose change respectively. In each row, the leftmost image is the input, and other images are the continuous manipulation results.

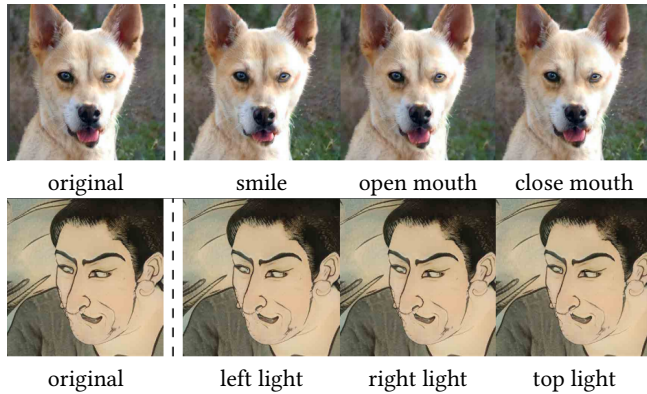


Fig. 12. Failure cases in the AFHQ Dog and Japan Ukiyo-e Face domain as for expression and illumination editing, respectively. This failure is because the datasets used for finetuning the StyleGAN2 do not contain enough variations for these attributes.

Nets. In *NeurIPS* 2014. 2672–2680.  
 Sungjoo Ha, Martin Kersner, Beomsu Kim, Seokjun Seo, and Dongyoung Kim. 2020. MarioNETT: Few-Shot Face Reenactment Preserving Identity of Unseen Targets. In *AAAI* 2020. 10893–10900.  
 Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. 2020. GANSpace: Discovering Interpretable GAN Controls. In *NeurIPS* 2020.  
 Jialu Huang, Jing Liao, and Sam Kwong. 2020. Unsupervised Image-to-Image Translation via Pre-trained StyleGAN2 Network. (2020). arXiv:2010.05713 [cs.CV]

Xun Huang and Serge J. Belongie. 2017. Arbitrary Style Transfer in Real-Time with Adaptive Instance Normalization. In *ICCV* 2017. 1510–1519.  
 Tero Karras, Samuli Laine, and Timo Aila. 2019. A Style-Based Generator Architecture for Generative Adversarial Networks. In *CVPR* 2019. 4401–4410.  
 Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and Improving the Image Quality of StyleGAN. In *CVPR* 2020. 8107–8116.  
 Hyeonwoo Kim, Mohamed Elgharib, Michael Zollhöfer, Hans-Peter Seidel, Thabo Beeler, Christian Richardt, and Christian Theobalt. 2019. Neural style-preserving visual dubbing. *ACM Trans. Graph.* 38, 6 (2019), 178:1–178:13.  
 Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Nießner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. 2018. Deep video portraits. *ACM Trans. Graph.* 37, 4 (2018), 163:1–163:14.  
 Ryota Natsume, Tatsuya Yatagawa, and Shigeo Morishima. 2018. FSNet: An Identity-Aware Generative Model for Image-Based Face Swapping. In *ACCV* 2018, Vol. 11366. 117–132.  
 Yuval Nirkin, Yosi Keller, and Tal Hassner. 2019. FSGAN: Subject Agnostic Face Swapping and Reenactment. In *ICCV* 2019. 7183–7192.  
 Justin N. M. Pinkney. 2020. Aligned Ukiyo-e faces dataset. <https://www.justinpinkney.com/ukiyoe-dataset>.  
 Justin N. M. Pinkney and Doron Adler. 2020. Resolution Dependent GAN Interpolation for Controllable Image Synthesis Between Domains. (2020). arXiv:2010.05334 [cs.CV]  
 Alec Radford, Luke Metz, and Soumith Chintala. 2016. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. In *ICLR* 2016.  
 Ravi Ramamoorthi and Pat Hanrahan. 2001. An efficient representation for irradiance environment maps. In *SIGGRAPH* 2001. 497–500.  
 Scott Schaefer, Travis McPhail, and Joe D. Warren. 2006. Image deformation using moving least squares. *ACM Trans. Graph.* 25, 3 (2006), 533–540.  
 Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. 2020. Interpreting the Latent Space of GANs for Semantic Face Editing. In *CVPR* 2020. 9240–9249.  
 Yujun Shen and Bolei Zhou. 2020. Closed-Form Factorization of Latent Semantics in GANs. *CoRR* abs/2007.06600 (2020).  
 Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019a. Animating Arbitrary Objects via Deep Motion Transfer. In *CVPR* 2019. 2377–2386.

- Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019b. First Order Motion Model for Image Animation. In *NeurIPS 2019*, 7135–7145.
- Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. 2017. Synthesizing Obama: learning lip sync from audio. *ACM Trans. Graph.* 36, 4 (2017), 95:1–95:13.
- Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. 2020a. StyleRig: Rigging StyleGAN for 3D Control Over Portrait Images. In *CVPR 2020*. 6141–6150.
- Ayush Tewari, Mohamed Elgharib, Mallikarjun B. R., Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. 2020b. PIE: portrait image embedding for semantic control. *ACM Trans. Graph.* 39, 6 (2020), 223:1–223:14.
- Ayush Tewari, Michael Zollhöfer, Hyeonwoo Kim, Pablo Garrido, Florian Bernard, Patrick Pérez, and Christian Theobalt. 2017. MoFA: Model-Based Deep Convolutional Face Autoencoder for Unsupervised Monocular Reconstruction. In *ICCV 2017*. 3735–3744.
- Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: image synthesis using neural textures. *ACM Trans. Graph.* 38, 4 (2019), 66:1–66:12.
- Justus Thies, Michael Zollhöfer, Matthias Nießner, Levi Valgaerts, Marc Stamminger, and Christian Theobalt. 2015. Real-time expression transfer for facial reenactment. *ACM Trans. Graph.* 34, 6 (2015), 183:1–183:14.
- Olivia Wiles, A. Sophia Koepke, and Andrew Zisserman. 2018. X2Face: A Network for Controlling Face Generation Using Images, Audio, and Pose Codes. In *ECCV 2018*, Vol. 11217. 690–706.
- Zongze Wu, Dani Lischinski, and Eli Shechtman. 2020. StyleSpace Analysis: Disentangled Controls for StyleGAN Image Generation. *CoRR abs/2011.12799* (2020).
- Sitao Xiang, Yuming Gu, Pengda Xiang, Mingming He, Koki Nagano, Haiwei Chen, and Hao Li. 2020. One-Shot Identity-Preserving Portrait Reenactment. *CoRR abs/2004.12452* (2020).
- Sicheng Xu, Jiaolong Yang, Dong Chen, Fang Wen, Yu Deng, Yunde Jia, and Xin Tong. 2020. Deep 3D Portrait From a Single Image. In *CVPR 2020*. 7707–7717.
- Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor S. Lempitsky. 2019. Few-Shot Adversarial Learning of Realistic Neural Talking Head Models. In *ICCV 2019*. 9458–9467.
- Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. 2020. In-Domain GAN Inversion for Real Image Editing. In *ECCV 2020*, Vol. 12362. 592–608.
- Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and Stan Z. Li. 2016. Face Alignment Across Large Poses: A 3D Solution. In *CVPR 2016*. 146–155.
- Michael Zollhöfer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Pérez, Marc Stamminger, Matthias Nießner, and Christian Theobalt. 2018. State of the Art on Monocular 3D Face Reconstruction, Tracking, and Applications. *Comput. Graph. Forum* 37, 2 (2018), 523–550.