

# DyStyle: Dynamic Neural Network for Multi-Attribute-Conditioned Style Editing

Bingchuan Li\*, Shaofei Cai\*, Wei Liu, Peng Zhang, Miao Hua, Qian He, Zili Yi †

ByteDance Ltd.

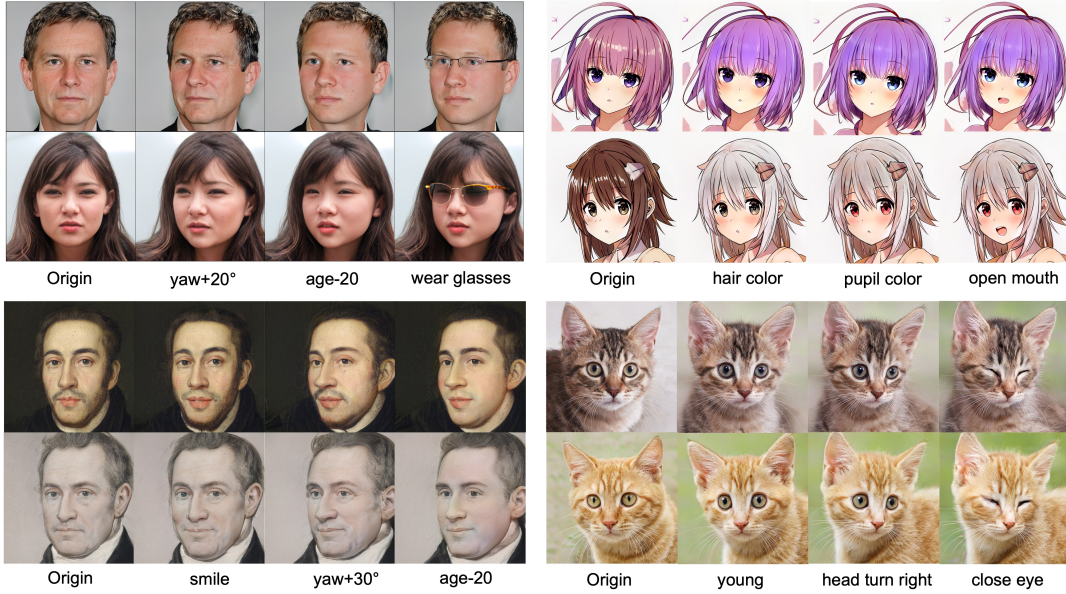


Figure 1: The multi-attribute-conditioned image editing results with our approach, achieved by manipulating the latent codes of four different pre-trained StyleGAN2 models.

## Abstract

Great diversity and photorealism have been achieved by unconditional GAN frameworks such as StyleGAN and its variations (Karras, Laine, and Aila 2019; Karras et al. 2020b,a). In the meantime, persistent efforts have been made to enhance the semantic controllability of StyleGANs. For example, a dozen of style manipulation methods have been recently proposed to perform attribute-conditioned style editing. Although some of these methods work well in manipulating the style codes along one attribute, the control accuracy when jointly manipulating multiple attributes tends to be problematic. To address these limitations, we propose a Dynamic Style Manipulation Network (DyStyle) whose structure and parameters vary by input samples, to perform nonlinear and adaptive manipulation of latent codes for flexible and precise attribute control. Additionally, a novel easy-to-hard training procedure is introduced for efficient and stable training of the DyStyle network. Extensive experiments have been con-

ducted on faces and other objects. As a result, our approach demonstrates fine-grained disentangled edits along multiple numeric and binary attributes. Qualitative and quantitative comparisons with existing style manipulation methods verify the superiority of our method in terms of the attribute control accuracy and identity preservation without compromising the photorealism. The advantage of our method is even more significant for joint multi-attribute control. The source codes are made publicly available at [phycvgan/DyStyle](https://github.com/byte-dance-research/phycvgan/DyStyle).

## Introduction

Recent development in Generative Adversarial Networks (GANs) (Goodfellow et al. 2014) has provided a new paradigm for realistic image generation. As one of the most celebrated GAN frameworks, StyleGAN (Karras, Laine, and Aila 2019) and the upgraded StyleGAN2 (Karras et al. 2020b,a, 2021), can produce diverse and high fidelity images with unmatched photorealism. Nevertheless, like other unconditional GAN frameworks, due to the highly-entangled nature of the latent space, StyleGAN offers users limited op-

\*C0-first author

†Corresponding author

tions to control over the attributes or semantics of generated images. Performing edits along one attribute can easily result in unwanted changes along other attributes.

To enhance the controllability of GAN-generated images, one stream of research is to manipulate the latent codes of unconditional GANs, without retraining the generator. Methods of this stream (Shen et al. 2020; Härkönen et al. 2020; Shen and Zhou 2020; Hou et al. 2020; Tewari et al. 2020; Abdal et al. 2020; Wang, Yu, and Fritz 2021; Xia et al. 2021; Roich et al. 2021; Alaluf, Patashnik, and Cohen-Or 2021b; Ren et al. 2021; Lang et al. 2021; Wu, Lischinski, and Shechtman 2021; Patashnik et al. 2021) attempt to achieve controlled image synthesis by exploring the semantics in the latent space of well-trained GANs. For example, some researchers worked on the hypothesis that the StyleGAN latent space is actually linear and they propose linear manipulations in that space. Three noteworthy works are InterFaceGAN (Shen et al. 2020), GANSpace (Härkönen et al. 2020) and SeFa (Shen and Zhou 2020). These methods either takes a data driven approach and uses Principal Component Analysis (PCA) to learn the most important directions. Specifically, GANSpace (Härkönen et al. 2020) identifies latent directions based on applied PCA either in latent space or feature space, and interpretable controls can be defined by layer-wise perturbation along the principal directions. Upon analyzing these directions they discover that the directions correspond to meaningful semantic edits. SeFa (Shen and Zhou 2020) is also an unsupervised linear method that directly decomposes the pre-trained weights rather than the latent space or feature space. Unlike GANSpace and SeFa, InterFaceGAN (Shen et al. 2020) use labeled data to learn a SVM to discover the separation plane and the directions of certain attributes. The assumption of a linear latent space is a useful simplification that produces fair results for limited types of attributes. However, we are still unable to produce precisely controlled results. Due to the entanglement between different semantics in the latent space, performing edits along one attribute could lead to unexpected changes of other semantics.

Seeing the limitation of linear decomposition of latent space, some researchers intend to process nonlinear edits of the latent codes (Tewari et al. 2020; Hou et al. 2020; Abdal et al. 2020; Wang, Yu, and Fritz 2021). The most notable works are StyleRig (Tewari et al. 2020), StyleFlow (Abdal et al. 2020) and Hijack-GAN (Wang, Yu, and Fritz 2021). Specifically, StyleRig (Tewari et al. 2020) employs a 3D morphable face models (3DMMs) (Blanz and Vetter 1999), and a rigging network to map 3DMM’s semantic parameters to StyleGAN’s input. In this way, they condition the edits of style codes on the 3DMM’s semantic parameters, which include information of illumination, expression, camera pose. The success of StyleRig relies on the massive paired data generated by manipulating the 3DMMs in 3D domain. Although StyleRig generates very nice results for the manipulation of head pose and illumination, the detailed control of other facial attributes ultimately did not work. StyleFlow (Abdal et al. 2020) seeks continuous and nonlinear normalizing flows in the latent space conditioned by attribute features. However, the training of StyleFlow does not explicitly

enforce the control accuracy of attributes, and the assumption of normal distribution with respect to any attribute configuration does not always hold due to data biases. In addition, StyleFlow attempts to edit the limited  $\mathbf{W}$  space rather than the broader  $\mathbf{W}+$  space. Hijack-GAN (Wang, Yu, and Fritz 2021) assumes that the style edits along different attribute should be orthogonal to each other and thoroughly explores the local orthogonality across the nonlinear multi-attribute space. However, the assumption of orthogonality relies heavily on the linear assumption of the style space, and as a result the control accuracy is very limited. Generally speaking, nonlinear methods exhibit superior semantic disentanglement and attribute control to linear methods. Nevertheless, the control accuracy and identity preservation of existing nonlinear methods is still below the standards of design industry. Moreover, jointly manipulating multiple attributes typically results in larger control error: see Figure 18.

We argue that the style editing network should be able to adapt to wide varieties of attribute configurations, and when training the style editing network for multiple attribute manipulation, any biases of the distribution of attribute configurations would easily result in systematic control errors. Take realistic portrait editing as an example, for one case we only change the hair color of the portrait. For another case, we change both the hair color and age of the portrait. We wish the two attribute configurations are sampled with equal probability and are both well handled by the style editing network. We made two efforts to address this issue, we evenly sample the attribute configurations during training instead of using static set of training samples. To make this possible, we employ pre-trained knowledge networks to provide “on-the-fly” supervisions rather than using static labels. In addition, we employ a dynamic style editing network consisting of multiple experts, each of which is responsible for the manipulation of one attribute. We dynamically activate a subset of the experts based on whether the corresponding attributes are edited or not. That means, the structure and parameters of the style editing network vary by different input samples. Experiments show that the Dynamic Style Manipulation Network (DyStyle) that processes each sample with data-dependent architectures and parameters can well adapt to various types of attribute configurations, though the training of the network is a bit tricky. To make the training of the DyStyle network easier, we adopt a novel easy-to-hard training procedure in which the DyStyle network is trained for editing a single attribute at a time, and then trained for jointly manipulating multiple randomly-sampled attributes. Generally speaking, our contributions include:

- A Dynamic Style Manipulation Network (DyStyle) is proposed to perform multi-attribute-conditioned editing of the StyleGAN latent codes, whose structure and parameters vary by input samples, leveraging its adaptability to wide varieties of attribute configurations.
- We train the DyStyle network by following a two-stage easy-to-hard training procedure, which proves to be beneficial to avoid getting trapped on local minimum and increase the training stability.



- Comprehensive evaluations on various datasets (realistic faces, comics, artistic portraits, animal faces) demonstrate improved attribute control accuracy and better identity preservation of our approach over existing static architectures. The improvements are more significant when jointly manipulating the styles along multiple attributes.

## Related Work

### Unconditional GANs

Generative Adversarial Network (GAN) is first introduced by Goodfellow et al. (Goodfellow et al. 2014), and has been one of the most active fields in deep neural networks. One research direction is to improve the GAN architectures, loss functions, and training dynamics for improved quality, diversity and stability of training. In terms of the GAN architectures, DCGAN (Radford, Metz, and Chintala 2015), ProgressiveGAN (Karras et al. 2017), StyleGAN (Karras, Laine, and Aila 2019), BigGAN (Brock, Donahue, and Simonyan 2018) and StyleGAN2 (Karras et al. 2020b) architectures are the top known architectures in history of development. StyleGAN (Karras, Laine, and Aila 2019) and its upgraded version StyleGAN2 (Karras et al. 2020b,a, 2021) are especially strong in synthesizing realistic faces. We build our work on StyleGAN2, as it achieves the best photorealism.

### Conditional GANs for controllable face synthesis

Conditional GANs (Mirza and Osindero 2014) have given rise to many image manipulation applications. Unlike unconditional GANs that take random noises as input, conditional GANs take meaningful priors (e.g., attribute, 3D model parameters, sketches, label maps) as input and synthesizes relevant images, thus offering users certain level of control.

In the context of faces, many conditional GANs (Yang et al. 2021; Kim et al. 2021; Chan et al. 2021; Chen et al. 2021; Zhu et al. 2020b; Shoshan et al. 2021; Li et al. 2021; Deng et al. 2020; Choi et al. 2020) are proposed to generate or manipulate a face based on attributes or other semantics. One of the most notable works is StarGAN (Choi et al. 2018, 2020) that proposes a GAN architecture that condition face generating on facial attributes such as hair color, gender, and age. Some recent works such GAN-control (Shoshan et al. 2021), HiSD (Li et al. 2021), ConfigNet (Kowalski et al. 2020), DiscoFaceGAN (Deng et al. 2020) continue to work in this direction, and achieve better attribute disentanglement and improved image quality. Scribber (Sangkloy et al. 2017), FaceShop (Portenier et al. 2018), SC-FEGAN (Jo and Park 2019) and DeepFillV2 (Yu et al. 2019) use generators conditioned on sketches and color information. Some other models focus on specialized image manipulation techniques, such as makeup transfer (Jiang et al. 2020), motion transfer (Siarohin et al. 2019; Wiles, Koepke, and Zisserman 2018; Zakharov et al. 2019; Wang, Mallya, and Liu 2020; Li et al. 2019; Siarohin et al. 2020; Wang et al. 2018, 2019), novel view synthesis (Huang et al. 2017; Chan et al. 2020;

Gao et al. 2020; Tan et al. 2020; Chan et al. 2020), relighting (Zhou et al. 2019; Egger et al. 2018; Sengupta et al. 2018; Zhou et al. 2019; Zhang et al. 2020), gaze manipulation (Ganin et al. 2016) or hair editing (Ak et al. 2019). Conditional GANs play an essential role in the task of controlled image generation and semantic editing. Whereas the resultant image qualities of conditional GANs still fail to match the realism produced by unconditional GANs like StyleGAN (Karras, Laine, and Aila 2019; Karras et al. 2020b).

### Image-to-style mapping

The techniques that embed Images into latent space of unconditional GANs facilitate various tasks including GAN-based image analysis (Donahue, Krähenbühl, and Darrell 2016; Donahue and Simonyan 2019; Mukherjee et al. 2019), image manipulation (Zhu et al. 2020a; Richardson et al. 2020; Abdal, Qin, and Wonka 2019, 2020; Tov et al. 2021), enhancement (Menon et al. 2020) and image compression (Agustsson et al. 2019). One straightforward way to do this is to train an encoder network that maps an image into the latent space. Another technique is to use an optimization algorithm to search the latent codes that can minimize the difference between the generated image and the target (Richardson et al. 2020; Abdal, Qin, and Wonka 2019, 2020). A third way is to combine the two techniques. That is, an encoder is first used to predict an approximate neighbor and the neighbor is set as the initial point of the optimization algorithm. In the context of StyleGAN (Choi et al. 2018) and StyleGAN2 (Choi et al. 2020), the optimization based technique does not work well for the extended latent space  $\mathbf{W}+$  and typically exerts larger reconstruction error. In this paper, we will use an encoder-based method (Richardson et al. 2020; Tov et al. 2021) to locate the latent code of a real photograph for further manipulation.

## Method

### Framework

As shown in Figure 2, our method manipulates the extended latent code  $\mathbf{W}+$  using a Dynamic Style Manipulation Network (DyStyle). The extended latent code  $\mathbf{W}+$  consists of 18 different 512-dimensional vectors, one for each input layer of StyleGAN2 generator (Karras et al. 2020b). Compared with  $\mathbf{W}$ ,  $\mathbf{W}+$  significantly expands the latent space and results in smaller reconstruction error when designated to represent real photographs (Richardson et al. 2020). The  $\mathbf{W}+$  can be either mapped from a random Gaussian noise vector  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{1})$  with the Style Mapping Network of StyleGAN2 (Karras et al. 2020b) or embedded from a real photograph with the image-to-style encoder of pSp (Richardson et al. 2020) or E4E (Tov et al. 2021). The DyStyle network takes an attribute specification and  $\mathbf{W}+$  as inputs and predicts a manipulated latent code  $\hat{\mathbf{W}}+$ . The attribute specification  $\mathbf{Attr}$  is made up of a set of attribute values specified by the user, defining the appearance of the desired image. The attribute set, in our experimental setting, includes numeric attributes (e.g., yaw, pitch of the head pose, age) and binary attributes (e.g., glasses, smile, black hair, mustache,

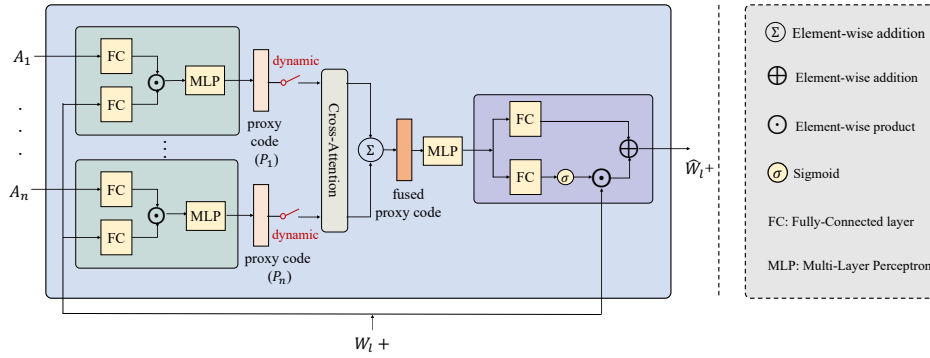


Figure 2: The framework and training losses of our multi-attribute-conditioned style editing approach. Note that the Encoder is a pre-trained image-to-style encoder provided by pSp (Richardson et al. 2020) or E4E (Tov et al. 2021). The Dynamic Network is trained for attribute-conditioned image editing while the StyleGAN2 Generator, the Style Mapping Network and the Encoder are held fixed.

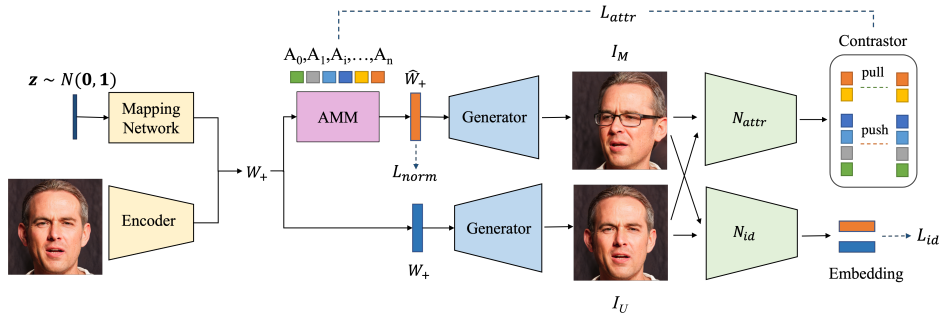


Figure 3: The architecture of our Dynamic Style Manipulation Network (DyStyle).  $\mathbf{W}_{l+}$  ( $l \in \{0, 1, \dots, 17\}$ ) is one of the 18 vectors of  $\mathbf{W}_+$ . Multiple experts are employed, in which each expert is responsible for the processing of one attribute before they are joined. Whether an expert is activated is based on whether the corresponding attribute is intended for editing or not.

close eye, open mouth). The number of attributes can be expanded without modifying the framework.

The manipulated latent code  $\hat{\mathbf{W}}_+$  is fed to the StyleGAN2 generator to generate the corresponding manipulated image  $\mathbf{I}_M$ . In the meantime, the original latent code  $\mathbf{W}_+$  is mapped to the untouched image  $\mathbf{I}_U$ .  $\mathbf{I}_M$  is expected to maintain the identity of  $\mathbf{I}_U$ , while matching the target attribute specification. Such constraint is enforced with the pre-trained attribute predictors  $N_{attr}$  and a pre-trained identity recognition model  $N_{id}$ .

### Dynamic network architecture

The architecture of DyStyle is shown in Figure 3. As shown, the DyStyle network manipulates each  $\mathbf{W}_{l+}$  code separately, by taking the attribute configuration and  $\mathbf{W}_{l+}$  as input and predicting the proxy code  $P$  which is further used to linearly modulate the  $\mathbf{W}_{l+}$  code itself. By conditioning the proxy code jointly upon the attribute configuration and  $\mathbf{W}_{l+}$ , rather than the attribute configuration solely, the network is able to predict the proxy code adaptively for each input case of  $\mathbf{W}_+$ , rather than generate a uniform modulation parameters for all cases of  $\mathbf{W}_+$ . The DyStyle network employs multiple experts to process the attributes separately before they fused with cross-attention and element-wise ad-

dition, so that each expert can focus on the processing of one attribute before inter-communication is performed. The use of cross-attention is for two considerations. First, it allows information communication across experts, which encourages better attribute disentanglement. For example, the edits of yaw and pitch may bring changes to the same regions (e.g. nose or mouth), thus resulting in correlation or entanglement between them. When jointly manipulating yaw and pitch, the edits of yaw should also take the edits of pitch into account. Cross-attention allows different experts to communicate to each other and enables an expert adapt to the unexpected influences caused by the edits of other attributes. This argument is well supported by the ablation studies: see the supplementary materials for more details. Second, the cross-attention module is well suited for a variable number of attributes, which is an important feature of our dynamic architecture.

The DyStyle activates only a subset of the attribute branches based on whether the corresponding attributes are intended for editing or not. The switch of each attribute branch controls whether the corresponding branch is activated: see Figure 3. For example, when editing realistic portrait, for one case we only change the hair color of the portrait. Then we only need to activate the branch for hair color

editing. For another case, we change both the hair color and age of the portrait. Thus we activate both branches for age and hair color. As an extreme case, if no branches are activated, the input style code is not edited. The dynamic feature of the style editing network is valid both in the training phase and inference phase. With such a design, we wish to leverage the adaptiveness of the network, so that all kinds of attribute configurations be well handled by the style editing network. We have a comprehensive ablation studies on this feature and verify its effective in practice.

We join the features after attribute-specific processing by cross-attention and element-wise addition. Such a design favors disentangled attribute editing and improved control precision. The cross attention is computed as Eq. 1.

$$P_i = \sum_j V_j \odot \frac{\exp(Q_i \cdot K_j)}{\sum_j \exp(Q_i \cdot K_j)} \quad (1)$$

where  $Q_i, K_i, V_i$   $i \in \{1, 2, \dots, n\}$  are the query, key and value vectors computed from the proxy code  $P_i$  with an FC layer respectively. As an extreme case, when  $n = 1$ , the output of the cross-attention  $P_i = V_i$ .

Some more architecture features are elaborated in the supplementary materials.

## Training method

**Training losses** The DyStyle network is trained with an objective consisting of 3 types of losses, which is defined as

$$L = \alpha_{attr} L_{attr} + \alpha_{id} L_{id} + \alpha_{norm} L_{norm} \quad (2)$$

where  $L_{attr}$  is the attribute loss for various attributes (e.g., pose, age, black hair, glasses, smile), which are used to enforce the consistency of target attributes and measured attributes of the manipulated image  $\mathbf{I}_M$ .  $L_{id}$  is the identity loss intended to preserve the identity of the original image, while  $L_{norm}$  is the normalization loss discouraging degradation of image quality.  $\alpha_{attr}$ ,  $\alpha_{id}$  and  $\alpha_{norm}$  are the corresponding coefficients for each loss term.

The form of  $L_{attr}$  that relies on a set of contrastors differs for numeric attributes and binary ones. We employs a contrastive form of loss for numeric attributes, which turns out to be more tolerant to systematic errors of the pre-trained attribute estimators, and it proves to have superior performance in terms of the accuracy of attribute control (Deng et al. 2020; Shoshan et al. 2021).

Specifically, the contrastive loss for numeric attribute  $A^i$  is defined as

$$L_{attr}^A = \max\left(\left|A_M^i - A_U^i - \Delta_{A^i}^{gt}\right| - T_{A^i}, 0\right) \quad (3)$$

where  $A_M^i, A_U^i$  are attribute value of Face  $\mathbf{I}_M, \mathbf{I}_U$  measured by a pre-trained attribute estimation network  $N_{A^i}$ .  $\Delta_{A^i}^{gt}$  is the ground-truth pose variation of  $\mathbf{I}_M$  to  $\mathbf{I}_U$ , specified at the input of DyStyle.  $T_{A^i}$  are constant thresholds, which is set to 3 for yaw and pitch, and 5 for age.

For binary attributes, the input attribute values are either 0 or 1, and they represent the status of the target attributes

(1 means ‘‘with’’ and 0 means ‘‘without’’). We employ a pre-trained multi-task network ( $N_{A^i}$ ) to predict all the binary attributes of an image, thus  $\mathbf{A}_M^i = N_{A^i}(\mathbf{I}_M)$ ,  $\mathbf{A}_U^i = N_{A^i}(\mathbf{I}_U)$ .

The attribute loss of the binary attributes is written as

$$L_{attr}^A = \begin{cases} 1 - \frac{\mathbf{emb}_M^A \cdot \mathbf{emb}_U^A}{\|\mathbf{emb}_M^A\| \cdot \|\mathbf{emb}_U^A\|}, & \text{if } A_{gt}^i == A^i \\ -\sum_{A^i} [A_{gt}^i \log A_M^i + (1 - A_{gt}^i) \log(1 - A_M^i)], & \text{otherwise} \end{cases} \quad (4)$$

where  $\mathbf{emb}_M^A$  (or  $\mathbf{emb}_U^A$ ) is the activation of the second last layer of the pre-trained multi-attribute predictor  $N_{A^i}$  given  $\mathbf{I}_M$  (or  $\mathbf{I}_U$ ) as input.  $A_{gt}^i$  is the ground-truth (or target) attribute specification. The similarity score of  $\mathbf{emb}_M^A$  and  $\mathbf{emb}_U^A$  is enforced to 1 when no edits are intended. Further, the cross-entropy loss is calculated separately for each binary attribute, and finally summed up.

The original face  $\mathbf{I}_U$  and the manipulated face  $\mathbf{I}_M$  are expected to have the same identity. Therefore, the identity loss are computed as

$$L_{id} = 1 - \frac{\mathbf{emb}_M^{id} \cdot \mathbf{emb}_U^{id}}{\|\mathbf{emb}_M^{id}\| \cdot \|\mathbf{emb}_U^{id}\|} \quad (5)$$

where  $\mathbf{emb}_M^{id}$  (or  $\mathbf{emb}_U^{id}$ ) is the feature embedding of the face in  $\mathbf{I}_M$  (or  $\mathbf{I}_U$ ), extracted with a pre-trained face recognition model. The cosine similarity score of  $\mathbf{emb}_M^{id}$  and  $\mathbf{emb}_U^{id}$  are expected to be approximate to 1.

When the editing targets are not realistic faces, the conventional face recognizer does not serve as an effective identity representation. In this cases (e.g., comics or animal faces), we employ LPIPS loss (Zhang et al. 2018) as the identity loss, which is written as

$$L_{id} = \sum_l \frac{1}{H_l \cdot W_l} \sum_{h_l, w_l} \alpha_l \odot \|y_M^l - y_U^l\|_2^2 \quad (6)$$

where  $h_l, w_l$  are the coordinates of the  $l$ -th feature map, and  $y_M^l$  (or  $y_U^l$ ) is the activations of the  $l$ -th layer of a pre-trained VGG network upon input image  $\mathbf{I}_M$  (or  $\mathbf{I}_U$ ).  $\alpha_l \in \mathbb{R}^{C_l}$  is channel-wise scale vector for  $l$ -th layer activations.

The normalization loss is defined as

$$L_{norm} = \sum_l \|\hat{\mathbf{W}}_l - \mathbf{W}_{avg}\| \quad (7)$$

where  $\mathbf{W}_{avg}$  is the statistic center of the  $\mathbf{W}$  space of the pre-trained StyleGAN2 generator (Karras et al. 2020b), and  $\hat{\mathbf{W}}_l$  is the manipulated style vector corresponding to the  $l$ -th layer. As discussed in (Richardson et al. 2020), being closer to  $\mathbf{W}_{avg}$  means higher expected quality of the generated image.

Practically, we choose  $\alpha_{yaw} = 0.05$ ,  $\alpha_{pitch} = 0.05$ ,  $\alpha_{age} = 0.02$ ,  $\alpha_{A^i} = 1.0$  (if  $A^i$  is binary attribute),  $\alpha_{id} = 1.0$  and  $\alpha_{norm} = 0.001$  in our experiments.

**Two-stage training strategy** The DyStyle is trained with randomly sampled  $\mathbf{W}+$  codes with the Style Mapping Network and evenly-sampled attribute configurations. The training is conducted by following a two-phase procedure. In this



first stage, the network is trained for single-attribute manipulation, by randomly choosing an attribute for editing and evenly sample the target attribute value. That means, only an expert (or branch) for the edited attribute is activated at a time in this phase. Note that the loss is not changed. This phase allows each expert to focus on one attribute at a time and get used to easy editing cases. In the second stage, the DyStyle network is trained to adapt to situations when multiple attributes are manipulated jointly. In this phase, a combination of attribute set are randomly sampled and set as the input of the DyStyle network, so that the experts learn to communicate with each other and adapt to more complex attribute configurations.

## Experimental Results

### Experiment setups

We experimented our style manipulation network on four pre-trained StyleGAN2 models, the one trained on FFHQ dataset (Karras, Laine, and Aila 2019) for realistic face generation, the one trained on MetFace (Karras et al. 2020a) for artistic face synthesis, that trained on AFHQ (Choi et al. 2020) for animal face synthesis and the one trained on comics dataset (Branwen 2015) for comic face generation. As the training of our style manipulation network does not require any external images, the training samples are basically randomly-sampled attribute-configuration-and- $\mathbf{W}+$  pairs, which are generated on the fly during training. Specifically, the  $\mathbf{W}+$  vectors are mapped from randomly-sampled gaussian noise vector  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$  (truncated with 0.7) with the style mapping network of StyleGAN2 (Karras et al. 2020b), and the attribute configurations are produced by evenly sampling the attribute values along each attribute space (e.g.,  $age \sim U[-30, 30]$ ,  $smile \sim U\{0, 1\}$ ). The detailed attribute settings in our experiments for different datasets are listed in Table 4 in the supplementary materials.

As shown in Figure 2, we employ a couple of pre-trained attribute prediction networks to supervise the training. Specifically, we employ the official pre-trained HopeNet model (Ruiz, Chong, and Rehg 2018) implemented in PyTorch for pose estimation. The age estimator used for training is the pre-trained age regressor implemented in PyTorch (Rothe, Timofte, and Van Gool 2015). We employ the official pre-trained CircularFace model (Huang et al. 2020) for identity embedding extraction of realistic/artistic faces. The official pre-trained VGG-19 model for comic/animal identity embedding extraction. The multi-task binary attribute classifier is a ResNet34 (He et al. 2016) trained by ourselves on the CelebA dataset (Liu et al. 2015) (for realistic/artistic faces), AFHQ (for animal faces), and comics dataset (Branwen 2015) (for comic face).

Our style manipulation network is implemented in PyTorch 1.6. It is trained with batch size of 8 on a single Tesla V100 GPU. It is optimized using Adam optimizer (Kingma and Ba 2014) with  $\beta_1 = 0.5$  and  $\beta_2 = 0.99$ , and the learning rate is fixed at  $10^{-4}$ . In all experiments, our model is trained for 50,000 steps for single-attribute manipulation (Stage I) and then 100,000 additional epochs for for multi-attribute

editing (Stage II).

Various attribute manipulation results of our approach are presented in Figure 1. Some more attribute-controlled image generation results for each datasets are presented in Figure 11, 14 (realistic faces), Figure 12 (artistic faces), Figure 15 (Comics) and Figure 16 (animal faces) in the supplementary materials. Some more high-resolution ( $1024 \times 1024$ ) realistic face editing results by our approach are presented in Figure 13 in the supplementary materials. With the image-to-style encoder provided in (Richardson et al. 2020), we also conducted attribute-conditioned editing of real photos and present the results in Figure 17 in the supplementary materials.

### Comparisons

To verify the effectiveness of the proposed method, we compared our approach with the state-of-the-art style manipulation methods including InterFaceGAN (Shen et al. 2020), StyleFlow (Abdal et al. 2020). To assure fair comparisons, these methods are designated to manipulate the latent space of the same pre-trained StyleGAN2 model which was trained on FFHQ dataset (Karras, Laine, and Aila 2019). As InterFaceGAN does not offer the manipulation of pitch, mustache and hair color. the comparisons only account for the editing of four attributes (yaw, age, glasses and smile). As the separation plane of InterFaceGAN only specify the directions of attribute control, we roughly estimated the physical meaning of scales based on a few labeled examples.

With the same set of 5000 attribute-configuration-and- $\mathbf{W}+$  pairs as test set, we conduct quantitative evaluations on these methods. Specifically, we evaluate the accuracy (or precision) of attribute control, preservation of identity and image quality respectively with well-defined metrics. To evaluate the precision of attribute control, we employ the pre-trained attribute predictors to predict the attribute labels of manipulated images and then compare them with target labels. For assure fairness, we employ a different set of attribute predictors that are excluded from those used for training. Specifically, the pose estimator is the official pre-trained model from (Yang et al. 2018), and the age estimator is officially provided by (Alaluf, Patashnik, and Cohen-Or 2021a). The glasses and smile classifier is a multi-task ResNet50 classifier (He et al. 2016) trained by ourselves with CelebA dataset (Liu et al. 2015). With the predicted labels, we compute the Mean Absolute Error (MAE) of yaw and age, and the classification accuracy of “glasses” and “smile” attributes. As for identity preservation, we calculate the average cosine similarity score of manipulated faces and original ones: see Table 1. In terms of image quality, we evaluate the distance between the distribution of manipulated images and that of real images (FFHQ dataset) with the Fréchet Inception Distance (FID) (Heusel et al. 2017); see Table 2. As shown, our model exhibits higher control precision of all attributes except for age. In addition, our approach achieves higher average identity similarity score in all contexts. When joint manipulating multiple attributes, the attribute control precision and identity similarity scores of InterFaceGAN (Shen et al. 2020) and StyleFlow (Abdal

Table 1: Quantitative comparisons of different attribute-conditioned style editing approached in terms of the identity preservation and attribute control accuracy. Note that we compute Mean Absolute Error for numeric attributes (e.g., yaw, age) and Classification Acc. for binary attributes (e.g., glasses, smile). The identity similarity score is between the original face and the manipulated.

attribute type	single-attribute editing				multi-attribute editing			
	yaw	age	glass	smile	glass+smile	yaw+glass		
method	Identity Similarity Score $\uparrow$							
InterFaceGAN (Shen et al.)	0.78 $\pm$ 0.05	0.82 $\pm$ 0.06	0.84 $\pm$ 0.12	0.95 $\pm$ 0.05	0.74 $\pm$ 0.05	0.68 $\pm$ 0.15		
StyleFlow (abdal et al.)	0.82 $\pm$ 0.07	0.86 $\pm$ 0.08	0.85 $\pm$ 0.1	0.96 $\pm$ 0.05	0.83 $\pm$ 0.12	0.78 $\pm$ 0.12		
<b>Ours</b>	<b>0.95 <math>\pm</math> 0.05</b>	<b>0.89 <math>\pm</math> 0.08</b>	<b>0.90 <math>\pm</math> 0.09</b>	<b>0.98 <math>\pm</math> 0.1</b>	<b>0.87 <math>\pm</math> 0.1</b>	<b>0.85 <math>\pm</math> 0.09</b>		
method	Attribute Control Accuracy							
	yaw $\downarrow$	age $\downarrow$	smile $\uparrow$	glass $\uparrow$	glass $\uparrow$	smile $\uparrow$	yaw $\downarrow$	glass $\uparrow$
InterFaceGAN (Shen et al.)	12.95	13.50	0.894	0.93	0.832	0.877	15.33	0.826
StyleFlow (Abdal et al.)	6.41	<b>12.78</b>	0.975	0.981	0.944	0.921	8.58	0.925
<b>Ours</b>	<b>6.33</b>	13.77	<b>0.976</b>	<b>0.988</b>	<b>0.963</b>	<b>0.955</b>	<b>7.26</b>	<b>0.961</b>

Table 2: The FID between manipulated faces  $I_M$  and real faces  $I_U$  from FFHQ dataset (Karras, Laine, and Aila 2019).

StyleGAN2	InterFaceGAN (Shen et al.)	StyleFlow (Abdal et al.)	<b>Ours</b>
26.89	61.6	52.23	<b>43.87</b>

et al. 2020) deteriorate significantly while our method performs consistently well. The FID of our method is slightly better than StyleFlow and InterFaceGAN, but worse than the vanilla StyleGAN2 probably due to the varying distribution of manipulated styles.

### Ablation studies

To verify the effectiveness of the proposed dynamic architecture and two-stage training procedure, we prepared two validation datasets separately for multi-attribute-conditioned realistic (FFHQ) editing and comic face editing. We started three trainings, which include the static architecture (all branches are activated regardless of the input as in Figure 3) trained for joint multi-attribute editing, the dynamic architecture trained with the two-stage training procedure, and that trained for multi-attribute editing only (single-stage). We visualize how the validation losses ( $L_{id}$ ,  $L_{attr}$ ,  $L_{attr}^A$  as defined in the method section) change against the number of training steps. As shown in Figure 4, for both experiments, the static architecture cannot converge as well as the dynamic architecture, implying its invulnerability in adapting to various kinds of attribute configurations. As for the dynamic architecture, the single-stage training procedure is highly unstable and achieves worse identity preservation and average control accuracy. Visual comparisons are illustrated in Figures 5, Figure 7 in the supplementary materials. More ablation studies on the architecture features and training techniques are presented in the supplementary materials.

### Conclusion

We propose a dynamic neural network that enables nonlinear and adaptive style manipulations for multi-attribute-conditioned image generation. As we only manipulate the latent space of StyleGAN2, our model is able to produce images with high quality. This is the advantage of our method

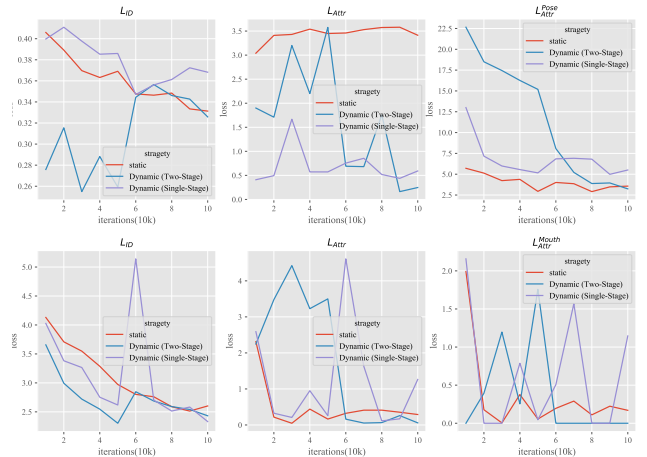


Figure 4: Comparisons of the dynamic architecture versus the static architecture, and the two-stage training versus single-stage training, experimented on FFHQ (top) and comic datasets (bottom). The validation losses demonstrate how well the model performs on the validation dataset in terms of the identity preservation and attribute control accuracy (the lower the better). As shown, as the two-stage training procedure trains for single-attribute editing in the first 50k iterations, its performance in terms of multi-attribute control accuracy at the beginning is not good. However, after 100k iterations, the two-stage training converges well and achieves the lowest  $L_{id}$  and  $L_{attr}$ .

over conditional GANs. Additionally, compared with other linear or nonlinear style manipulation approaches such as InterFaceGAN(Shen et al. 2020) and StyleFlow (Abdal et al. 2020), our model exhibits higher average precision of attribute-control and improved competency of identity preservation. When manipulating multiple attributes, the superiority of our approach becomes more significant.

However, even though our approach achieves high success rate of binary attribute editing, we still see shortcomings of the control precision of numeric attribute. For example, the MAE of age and pose are 12.77 (years) and 6.33 $^\circ$

respectively for realistic faces. Further, the average identity similarity score (about 0.85) as for multi-attribute editing can be further improved. In the future, we will persistently contribute to these directions.

## References

- Abdal, R.; Qin, Y.; and Wonka, P. 2019. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 4432–4441.
- Abdal, R.; Qin, Y.; and Wonka, P. 2020. Image2stylegan++: How to edit the embedded images? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8296–8305.
- Abdal, R.; Zhu, P.; Mitra, N.; and Wonka, P. 2020. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *arXiv e-prints*, arXiv:2008.
- Agustsson, E.; Tschannen, M.; Mentzer, F.; Timofte, R.; and Gool, L. V. 2019. Generative adversarial networks for extreme learned image compression. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 221–231.
- Ak, K. E.; Lim, J. H.; Tham, J. Y.; and Kassim, A. A. 2019. Attribute manipulation generative adversarial networks for fashion images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 10541–10550.
- Alaluf, Y.; Patashnik, O.; and Cohen-Or, D. 2021a. Only a Matter of Style: Age Transformation Using a Style-Based Regression Model. *arXiv preprint arXiv:2102.02754*.
- Alaluf, Y.; Patashnik, O.; and Cohen-Or, D. 2021b. ReStyle: A Residual-Based StyleGAN Encoder via Iterative Refinement. *arXiv preprint arXiv:2104.02699*.
- Blanz, V.; and Vetter, T. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 187–194.
- Branwen, G. 2015. Danbooru2019: A Large-Scale Crowd-sourced and Tagged Anime Illustration Dataset.
- Brock, A.; Donahue, J.; and Simonyan, K. 2018. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*.
- Chan, E. R.; Monteiro, M.; Kellnhofer, P.; Wu, J.; and Wetzstein, G. 2020. pi-GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis. *arXiv preprint arXiv:2012.00926*.
- Chan, E. R.; Monteiro, M.; Kellnhofer, P.; Wu, J.; and Wetzstein, G. 2021. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5799–5809.
- Chen, A.; Liu, R.; Xie, L.; Chen, Z.; Su, H.; and Yu, J. 2021. SofGAN: A Portrait Image Generator with Dynamic Styling. *ACM transactions on graphics*.
- Choi, Y.; Choi, M.; Kim, M.; Ha, J.-W.; Kim, S.; and Choo, J. 2018. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 8789–8797.
- Choi, Y.; Uh, Y.; Yoo, J.; and Ha, J.-W. 2020. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8188–8197.
- Deng, Y.; Yang, J.; Chen, D.; Wen, F.; and Tong, X. 2020. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5154–5163.
- Donahue, J.; Krähenbühl, P.; and Darrell, T. 2016. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*.
- Donahue, J.; and Simonyan, K. 2019. Large scale adversarial representation learning. *arXiv preprint arXiv:1907.02544*.
- Egger, B.; Schönborn, S.; Schneider, A.; Kortylewski, A.; Morel-Forster, A.; Blumer, C.; and Vetter, T. 2018. Occlusion-aware 3d morphable models and an illumination prior for face image analysis. *International Journal of Computer Vision*, 126(12): 1269–1287.
- Ganin, Y.; Kononenko, D.; Sungatullina, D.; and Lempit-sky, V. 2016. Deepwarp: Photorealistic image resynthesis for gaze manipulation. In *European conference on computer vision*, 311–326. Springer.
- Gao, C.; Shih, Y.; Lai, W.-S.; Liang, C.-K.; and Huang, J.-B. 2020. Portrait Neural Radiance Fields from a Single Image. *arXiv preprint arXiv:2012.05903*.
- Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*.
- Härkönen, E.; Hertzmann, A.; Lehtinen, J.; and Paris, S. 2020. Ganspace: Discovering interpretable gan controls. *arXiv preprint arXiv:2004.02546*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; and Hochreiter, S. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*.
- Hou, X.; Zhang, X.; Shen, L.; Lai, Z.; and Wan, J. 2020. GuidedStyle: Attribute Knowledge Guided Style Manipulation for Semantic Face Editing. *arXiv preprint arXiv:2012.11856*.
- Huang, R.; Zhang, S.; Li, T.; and He, R. 2017. Beyond face rotation: Global and local perception gan for photorealistic and identity preserving frontal view synthesis. In *Proceedings of the IEEE international conference on computer vision*, 2439–2448.
- Huang, Y.; Wang, Y.; Tai, Y.; Liu, X.; Shen, P.; Li, S.; Li, J.; and Huang, F. 2020. Curricularface: adaptive curriculum learning loss for deep face recognition. In *proceedings of*



- the *IEEE/CVF conference on computer vision and pattern recognition*, 5901–5910.
- Jiang, W.; Liu, S.; Gao, C.; Cao, J.; He, R.; Feng, J.; and Yan, S. 2020. Psgan: Pose and expression robust spatial-aware gan for customizable makeup transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5194–5202.
- Jo, Y.; and Park, J. 2019. SC-FEGAN: face editing generative adversarial network with user’s sketch and color. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 1745–1753.
- Karras, T.; Aila, T.; Laine, S.; and Lehtinen, J. 2017. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*.
- Karras, T.; Aittala, M.; Hellsten, J.; Laine, S.; Lehtinen, J.; and Aila, T. 2020a. Training generative adversarial networks with limited data. *arXiv preprint arXiv:2006.06676*.
- Karras, T.; Aittala, M.; Laine, S.; Härkönen, E.; Hellsten, J.; Lehtinen, J.; and Aila, T. 2021. Alias-Free Generative Adversarial Networks. *arXiv preprint arXiv:2106.12423*.
- Karras, T.; Laine, S.; and Aila, T. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4401–4410.
- Karras, T.; Laine, S.; Aittala, M.; Hellsten, J.; Lehtinen, J.; and Aila, T. 2020b. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8110–8119.
- Kim, H.; Choi, Y.; Kim, J.; Yoo, S.; and Uh, Y. 2021. StyleMapGAN: Exploiting Spatial Dimensions of Latent in GAN for Real-time Image Editing. *arXiv preprint arXiv:2104.14754*.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kowalski, M.; Garbin, S. J.; Estellers, V.; Baltrušaitis, T.; Johnson, M.; and Shotton, J. 2020. Config: Controllable neural face image generation. *arXiv preprint arXiv:2005.02671*.
- Lang, O.; Gandelsman, Y.; Yarom, M.; Wald, Y.; Elidan, G.; Hassidim, A.; Freeman, W. T.; Isola, P.; Globerson, A.; Irani, M.; et al. 2021. Explaining in Style: Training a GAN to explain a classifier in StyleSpace. *arXiv preprint arXiv:2104.13369*.
- Li, L.; Bao, J.; Yang, H.; Chen, D.; and Wen, F. 2019. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457*.
- Li, X.; Zhang, S.; Hu, J.; Cao, L.; Hong, X.; Mao, X.; Huang, F.; Wu, Y.; and Ji, R. 2021. Image-to-image Translation via Hierarchical Style Disentanglement. *arXiv preprint arXiv:2103.01456*.
- Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- Menon, S.; Damian, A.; Hu, S.; Ravi, N.; and Rudin, C. 2020. PULSE: Self-supervised photo upsampling via latent space exploration of generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2437–2445.
- Mirza, M.; and Osindero, S. 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
- Mukherjee, S.; Asnani, H.; Lin, E.; and Kannan, S. 2019. Clustergan: Latent space clustering in generative adversarial networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, 4610–4617.
- Patashnik, O.; Wu, Z.; Shechtman, E.; Cohen-Or, D.; and Lischinski, D. 2021. Styleclip: Text-driven manipulation of stylegan imagery. *arXiv preprint arXiv:2103.17249*.
- Portenier, T.; Hu, Q.; Szabo, A.; Bigdeli, S. A.; Favaro, P.; and Zwicker, M. 2018. Faceshop: Deep sketch-based face image editing. *arXiv preprint arXiv:1804.08972*.
- Radford, A.; Metz, L.; and Chintala, S. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
- Ren, X.; Yang, T.; Wang, Y.; and Zeng, W. 2021. Do Generative Models Know Disentanglement? Contrastive Learning is All You Need. *arXiv preprint arXiv:2102.10543*.
- Richardson, E.; Alaluf, Y.; Patashnik, O.; Nitzan, Y.; Azar, Y.; Shapiro, S.; and Cohen-Or, D. 2020. Encoding in style: a stylegan encoder for image-to-image translation. *arXiv preprint arXiv:2008.00951*.
- Roich, D.; Mokady, R.; Bermano, A. H.; and Cohen-Or, D. 2021. Pivotal Tuning for Latent-based Editing of Real Images. *arXiv preprint arXiv:2106.05744*.
- Rothe, R.; Timofte, R.; and Van Gool, L. 2015. Dex: Deep expectation of apparent age from a single image. In *Proceedings of the IEEE international conference on computer vision workshops*, 10–15.
- Ruiz, N.; Chong, E.; and Rehg, J. M. 2018. Fine-grained head pose estimation without keypoints. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2074–2083.
- Sangkloy, P.; Lu, J.; Fang, C.; Yu, F.; and Hays, J. 2017. Scribbler: Controlling deep image synthesis with sketch and color. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5400–5409.
- Sengupta, S.; Kanazawa, A.; Castillo, C. D.; and Jacobs, D. W. 2018. SfsNet: Learning Shape, Reflectance and Illuminance of Faces in the Wild’. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 6296–6305.
- Shen, Y.; Yang, C.; Tang, X.; and Zhou, B. 2020. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Shen, Y.; and Zhou, B. 2020. Closed-form factorization of latent semantics in gans. *arXiv preprint arXiv:2007.06600*.
- Shoshan, A.; Bhonker, N.; Kviatkovsky, I.; and Medioni, G. 2021. GAN-Control: Explicitly Controllable GANs. *arXiv preprint arXiv:2101.02477*.

- Siarohin, A.; Lathuilière, S.; Tulyakov, S.; Ricci, E.; and Sebe, N. 2019. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2377–2386.
- Siarohin, A.; Lathuilière, S.; Tulyakov, S.; Ricci, E.; and Sebe, N. 2020. First order motion model for image animation. *arXiv preprint arXiv:2003.00196*.
- Tan, Z.; Chai, M.; Chen, D.; Liao, J.; Chu, Q.; Yuan, L.; Tulyakov, S.; and Yu, N. 2020. MichiGAN: multi-input-conditioned hair image generation for portrait editing. *arXiv preprint arXiv:2010.16417*.
- Tewari, A.; Elgharib, M.; Bharaj, G.; Bernard, F.; Seidel, H.-P.; Pérez, P.; Zollhofer, M.; and Theobalt, C. 2020. Stylerig: Rigging stylegan for 3d control over portrait images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6142–6151.
- Tov, O.; Alaluf, Y.; Nitzan, Y.; Patashnik, O.; and Cohen-Or, D. 2021. Designing an encoder for stylegan image manipulation. *ACM Transactions on Graphics (TOG)*, 40(4): 1–14.
- Wang, H.-P.; Yu, N.; and Fritz, M. 2021. Hijack-gan: Unintended-use of pretrained, black-box gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7872–7881.
- Wang, T.-C.; Liu, M.-Y.; Tao, A.; Liu, G.; Kautz, J.; and Catanzaro, B. 2019. Few-shot video-to-video synthesis. *arXiv preprint arXiv:1910.12713*.
- Wang, T.-C.; Liu, M.-Y.; Zhu, J.-Y.; Liu, G.; Tao, A.; Kautz, J.; and Catanzaro, B. 2018. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*.
- Wang, T.-C.; Mallya, A.; and Liu, M.-Y. 2020. One-Shot Free-View Neural Talking-Head Synthesis for Video Conferencing. *arXiv preprint arXiv:2011.15126*.
- Wiles, O.; Koepke, A.; and Zisserman, A. 2018. X2face: A network for controlling face generation using images, audio, and pose codes. In *Proceedings of the European conference on computer vision (ECCV)*, 670–686.
- Wu, Z.; Lischinski, D.; and Shechtman, E. 2021. Stylespace analysis: Disentangled controls for stylegan image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12863–12872.
- Xia, W.; Yang, Y.; Xue, J.-H.; and Wu, B. 2021. TediGAN: Text-Guided Diverse Face Image Generation and Manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2256–2265.
- Yang, G.; Fei, N.; Ding, M.; Liu, G.; Lu, Z.; and Xiang, T. 2021. L2M-GAN: Learning To Manipulate Latent Space Semantics for Facial Attribute Editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2951–2960.
- Yang, T.-Y.; Huang, Y.-H.; Lin, Y.-Y.; Hsiu, P.-C.; and Chuang, Y.-Y. 2018. SSR-Net: A Compact Soft Stagewise Regression Network for Age Estimation. In *IJCAI*, volume 5, 7.
- Yu, J.; Lin, Z.; Yang, J.; Shen, X.; Lu, X.; and Huang, T. S. 2019. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 4471–4480.
- Zakharov, E.; Shysheya, A.; Burkov, E.; and Lempitsky, V. 2019. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9459–9468.
- Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 586–595.
- Zhang, X.; Barron, J. T.; Tsai, Y.-T.; Pandey, R.; Zhang, X.; Ng, R.; and Jacobs, D. E. 2020. Portrait shadow manipulation. *ACM Transactions on Graphics (TOG)*, 39(4): 78–1.
- Zhou, H.; Hadap, S.; Sunkavalli, K.; and Jacobs, D. W. 2019. Deep single-image portrait relighting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7194–7202.
- Zhu, J.; Shen, Y.; Zhao, D.; and Zhou, B. 2020a. In-domain gan inversion for real image editing. In *European Conference on Computer Vision*, 592–608. Springer.
- Zhu, P.; Abdal, R.; Qin, Y.; and Wonka, P. 2020b. Sean: Image synthesis with semantic region-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5104–5113.

## Appendix

### Architecture features

More features of the proposed DyStyle architecture are introduced in this section.

**Relative numeric attribute** In our design, the numeric attribute (e.g., age, yaw, pitch) values on which the style manipulation network is conditioned represent the relative change of the attribute with respect to the original image, rather than the absolute attribute values. For example, if the user specified  $yaw = +15$ , it means increasing the yaw angle by  $15^\circ$  with respect to the original face (i.e.,  $\Delta yaw = +15$ ). Such relative feature for numeric attributes has proved to be more effective in encouraging precise edits.

**Layer selection strategy** We adopt the idea used in (Abdal, Qin, and Wonka 2019, 2020; Abdal et al. 2020) to manipulate a subset of latent vectors for each attribute. This consideration has been justified by the observation that a specific attribute (e.g., head pose, age, glasses) are only correlated to a subset of latent vectors (Abdal, Qin, and Wonka 2019, 2020; Abdal et al. 2020). The layer selection mechanism is defined with the attribute-layer correlation factor  $\sigma_{l,A_i}$  ( $l \in \{0, 1, \dots, 17\}$ ,  $A_i \in \{\text{yaw, pitch, age, black-hair, glasses, smile, close-eye, open-mouth, ...}\}$ ). Specifically,  $\sigma_{l,x}$  are binary constants, which have the value of either 0 or 1. For example, yaw is only affected by three latent vectors  $\mathbf{W}_{0+}, \mathbf{W}_{1+}, \mathbf{W}_{2+}$ . Therefore, no proxy codes of  $\mathbf{W}_{l+}$  ( $i > 2$ ) are generated as for attribute yaw. The complete list of  $\sigma_{l,A_i}$  values is specified in Table 3.

### Extended ablations studies

Other than the dynamic architecture and two-stage training procedure, we conducted more ablation studies on other architecture features and training techniques used in our approach, and present the comparative results in the supplementary material. Specifically, the way of multi-attribute fusion (MLP or cross-attention) is studied and compared as in Figure 6. Figure 8 demonstrates if the relative attribute setting (only for numeric attributes) and the use of the contrastive losses is superior to the absolute attribute setting. Figure 10 examines how our model benefits from the identity loss  $L_{id}$  and normalization loss  $L_{norm}$  as introduced in the method section. Figure 9 compares two variations of DyStyle architectures. In our architecture design, we condition the proxy codes on both the latent code  $\mathbf{W}_{l+}$  and attribute specification to allow adaptive style modulation. Whereas, an alternative way is to condition the generation of proxy codes merely on the attributes. Figure 9 verifies the superiority of this feature. The visual comparisons of static architecture and dynamic architecture, two-stage training strategy and single-stage training strategy are presented in Figures 5, 7.

### Attribute settings and layer selection specifications.

The attribute settings for each datasets in our experiments are listed in Table 4. The full list of  $\sigma_{l,A_i}$  values with respect to layer order  $l$  and attribute type  $A_i$  is illustrated in Table 3.

## Visual experimental results

We present more attribute-controlled image generation results in Figure 11 (realistic faces), Figure 12 (artistic faces), Figure 15 (comic faces) and Figure 16 (animal faces). Some more high-resolution ( $1024 \times 1024$ ) realistic face editing results by our approach are presented in Figure 13, in the context of single-attribute manipulation and multi-attribute manipulation. Some more high-resolution ( $1024 \times 1024$ ) realistic face editing results by our approach are presented in Figure 13. The results of expression editing on realistic faces are illustrated in Figure 14.

With the image-to-style encoder provided in pSp (Richardson et al. 2020), we also conducted attribute-conditioned editing of real photos and present the results in Figure 17.

### Visual Comparisons between DyStyle and prior methods

We visually show some test results and demonstrate how the generated images vary by methods. As shown in Figure 18, when jointly manipulating multiple attribute, the preservation of identity and control accuracy along each attribute of prior methods are problematic. Some single-attribute editing results are demonstrated in Figure 19. As shown in Figure 19, StyleFlow (Abdal et al. 2020) exhibits good attribute disentanglement and identity preservation as ours. Whereas, the control precision of yaw and the smoothness of change when controlling binary attributes (glasses and smile) is inferior to ours. As InterFaceGAN (Shen et al. 2020) performs linear editing of style codes, unwanted change of identity and other attributes are noticeable. As all images are generated with the same StyleGAN2 generator, the degradation of image qualities is unnoticeable. Generally speaking, the performance gap in terms of single-attribute editing between existing methods and ours are not that noticeable as joint multi-attribute editing.

The comparisons of our method with other style editing method like StyleCLIP (Patashnik et al. 2021) is shown in Figure 20.



Table 3: Layer selection configurations: the value of  $\sigma_{l,A_i}$  with respect to layer order  $l$  and attribute type  $A_i$ .

$A_i$ $l$	realistic/artistic face							comic face				animal face				
	yaw	pitch	age	glasses	black hair	mustache	expressions	hair colors	pupil colors	open mouth	hair styles	head poses	young	breed	open mouth	close eye
0-1	1	1	1	0	0	0	0	0	0	0	0	1	1	0	0	1
2	1	1	1	1	0	0	0	0	0	0	0	1	1	0	0	1
3	0	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
4	0	0	1	0	0	0	1	0	0	1	0	0	1	0	0	1
5	0	0	1	0	0	0	1	0	0	1	0	0	1	1	0	1
6	0	0	1	0	0	0	0	0	0	1	1	0	1	1	0	1
7	0	0	1	0	0	1	0	0	0	1	1	0	1	0	1	1
8	0	0	0	0	1	1	0	0	0	1	1	0	1	0	1	1
9	0	0	0	0	1	1	0	0	0	0	1	0	0	0	0	0
10	0	0	0	0	0	0	0	1	1	0	1	0	0	0	0	0
11-14	0	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0
15-17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Table 4: The attribute settings in our experiments for different datasets.

	attribute name	type	value	detailed explanations
realistic/artistic face	yaw	numeric	$(-30,30)$	relative yaw change. "+20" means "increase yaw angle by 20°"
	pitch	numeric	$(-30,30)$	relative pitch change. "+20" means "increase pitch angle by 20°"
	age	numeric	$(-30,30)$	relative age change. "+20" means "become 20 years older"
	black-hair	binary	$\{0,1\}$	1 means having black hair
	mustache	binary	$\{0,1\}$	1 means having mustache
	expressions	Multi-class	$\{0,1\}^7$	(smile, angry, disgust, fear, sad, surprise, neutral) these expressions are exclusive. 1 means having that expressions.
	glasses	binary	$\{0,1\}$	1 means having glasses
comic face	pupil color	Multi-class binary	$\{0,1\}^8$	(red, yellow, blue, green, brown, purple, black, white) these colors are exclusive. 1 means the pupil is of that color.
	hair color	Multi-class binary	$\{0,1\}^8$	(red, yellow, blue, green, brown, purple, black, white) these colors are exclusive. 1 means the hair is of that color.
	open mouth	binary	$\{0,1\}$	1 means open mouth
	hair style	Multi-class binary	$\{0,1\}^2$	(long, short). the hair styles are exclusive. 1 means the hair is of that style.
animal face	head pose	Multi-class numeric	$\{0,1,2\}$	$\{0: \text{head turn left}, 1: \text{head facing front}, 2: \text{head turn right}\}$
	young	binary	$\{0,1\}$	0 mens young, 1 means old
	open mouth	binary	$\{0,1\}$	0 means open mouth. 1 means shut off the mouth.
	close eye	binary	$\{0,1\}$	1 means close eye
	breed	Multi-class binary	$\{0,1\}^5$	the breed set vary by dog or cat. Breed types are exclusive. 1 means the cat (or dog) is of that breed.

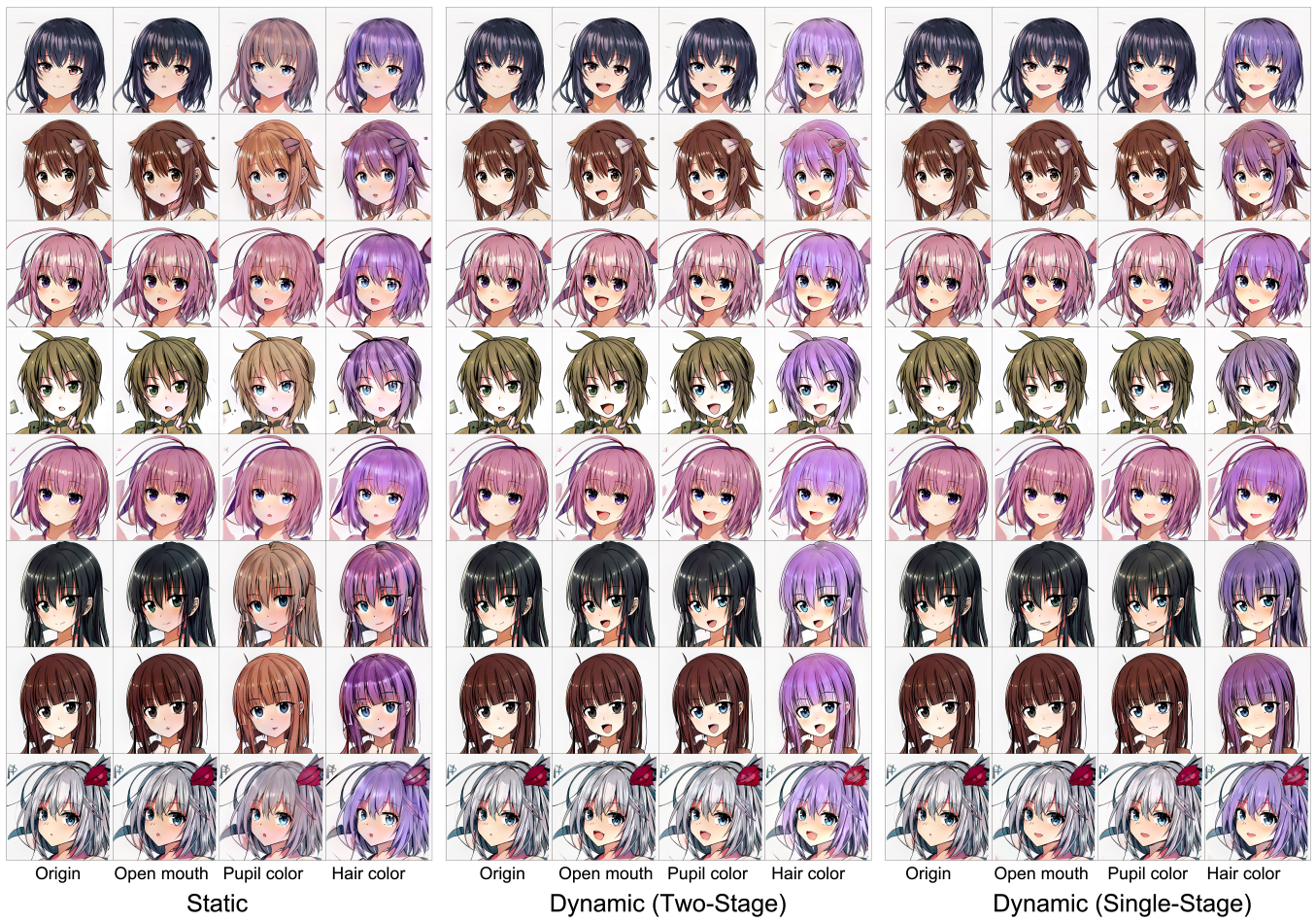


Figure 5: Qualitative comparisons of the dynamic architecture (**middle**) versus the static architecture (**left**), and the two-stage training (**middle**) versus single-stage training (**right**).



Figure 6: Qualitative comparisons of different multi-attribute fusion techniques (MLP versus Cross-Attention). Here MLP refers to that the cross-attention is removed and each branch is appended with an MLP before they are joined with element-wise addition. As shown, the severe entanglement of “glasses” and “yaw” can be observed from the results of MLP. Actually, we randomly tested 50 cases and observe that the entanglement between “yaw” and “glasses” occurs more frequently as for MLP than cross-attention (12/50 versus 3/50).



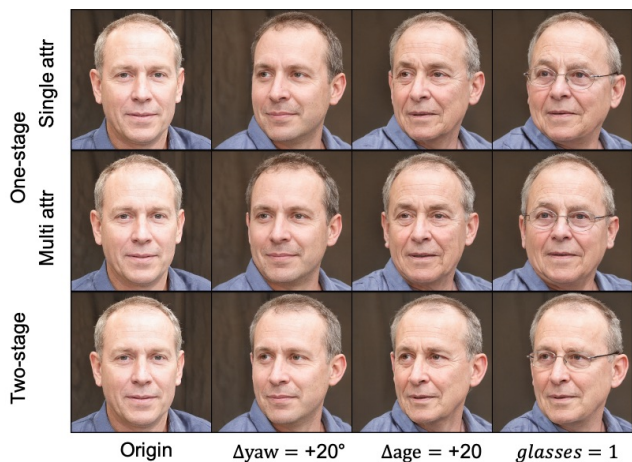


Figure 7: Ablation studies on the training strategy. We compare the results of two-stage training (single-attribute in Stage I and multi-attribute in Stage II, **(bottom)**) and one-stage training (single-attribute only **(top)**, multi-attribute only **(middle)**). As shown, the identity does not hold after multi-attribute editing **(top)** and the control of yaw is imprecise **(middle)**. The two-stage training strategy results in better editing results than one-stage training **(bottom)**.

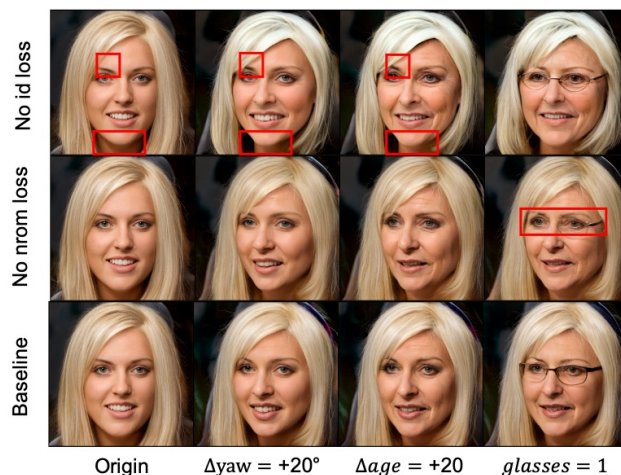


Figure 9: Ablation studies on the loss configuration. By removing the identity loss term or the normalization loss term from the full loss as in Eq 2, we retrain our model with the same hyper-parameters. Without  $L_{id}$  loss, the identity variation tends to be more significant **(top)**. Without  $L_{norm}$ , the generated images are prone to fall into failure modes **(middle)**: see the regions highlighted with red boxes.

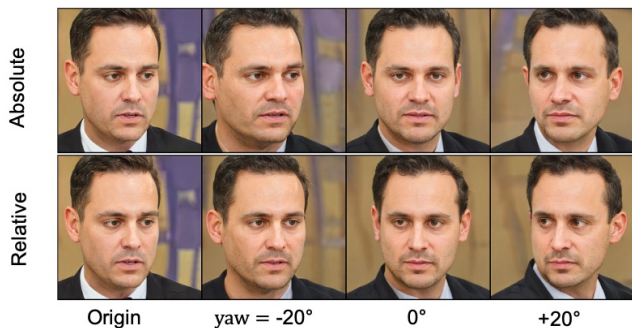


Figure 8: Comparisons of the relative attribute setting and the absolute setting. As shown, the absolute attribute setting results in unpleasant identity variation and imprecise control of head rotation along yaw.

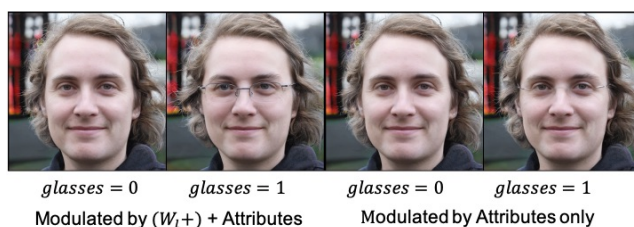


Figure 10: Ablation studies on the architecture design of the DyStyle. We compared the architecture conditioned the generation of proxy codes on the latent code and attributes (left), and that conditioned on attributes only (right). When changing the face from “no-glasses” to “with-glasses”, the left model generates faithful attribute editing results while the right model is prone to fail.

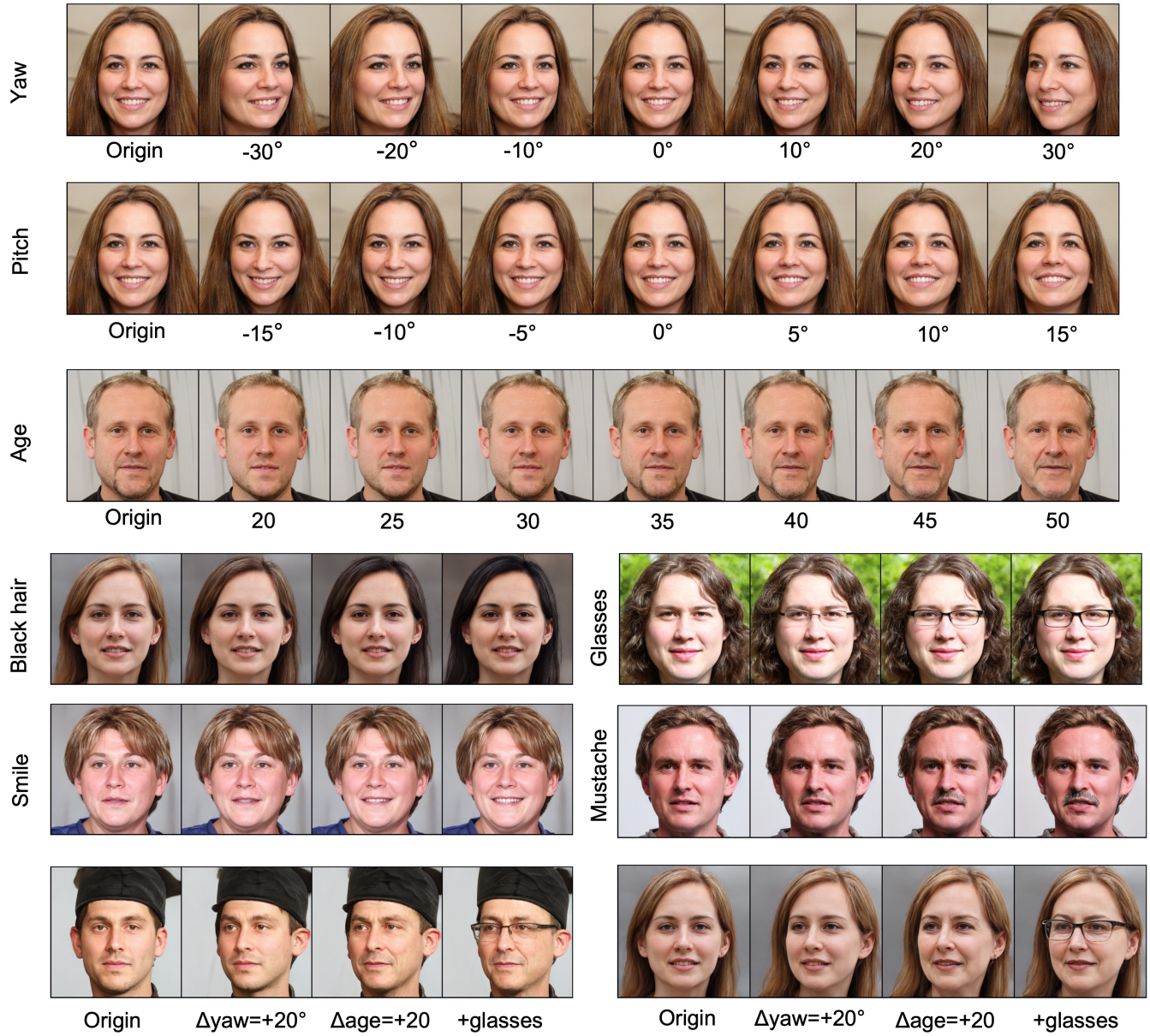


Figure 11: Results of single- and multi-attribute manipulation on realistic faces.



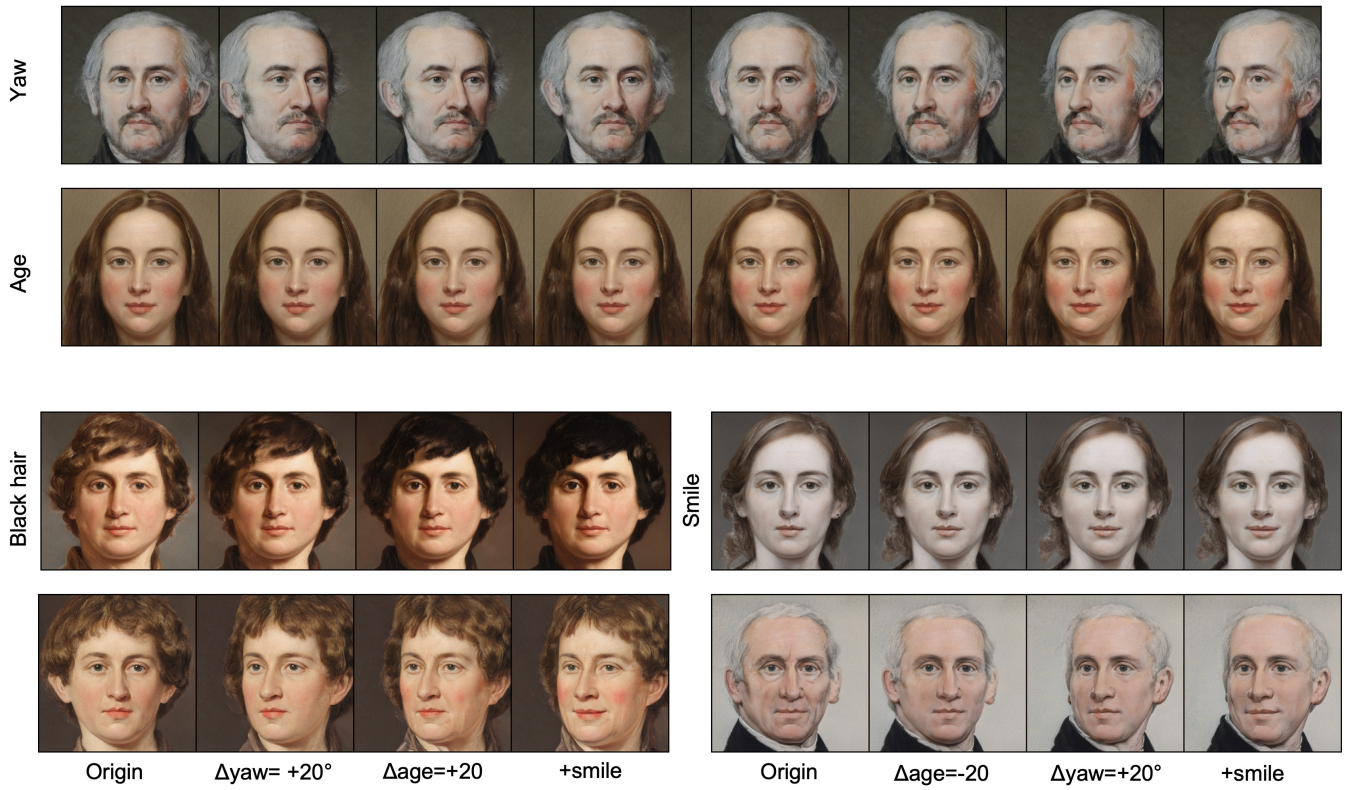


Figure 12: Results of single- and multi-attribute manipulation on artistic faces.



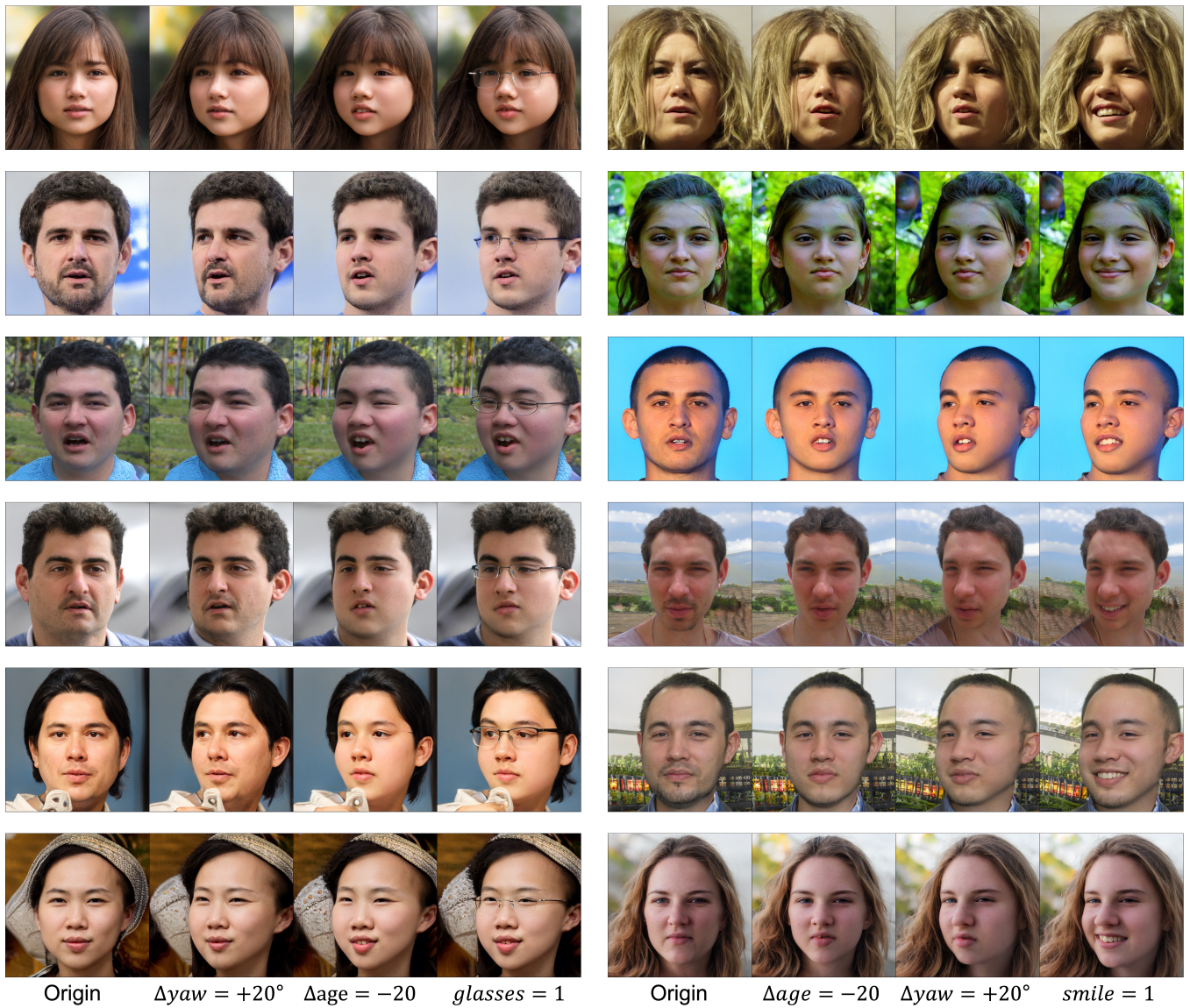


Figure 13: Results of multi-attribute manipulation on high-definition realistic faces (1024×1024).

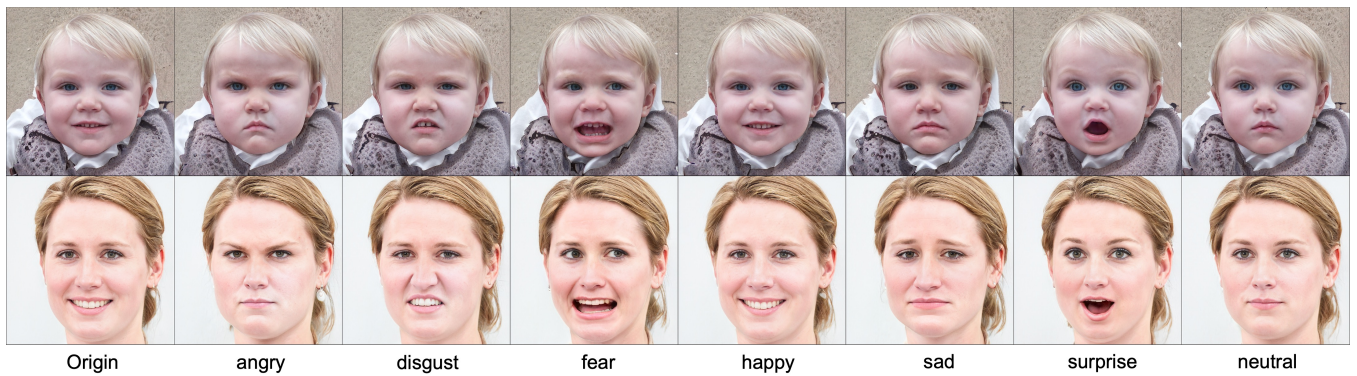


Figure 14: Results of multi-expression manipulation on realistic faces.





Figure 15: Results of single- and multi-attribute manipulation on comic faces.

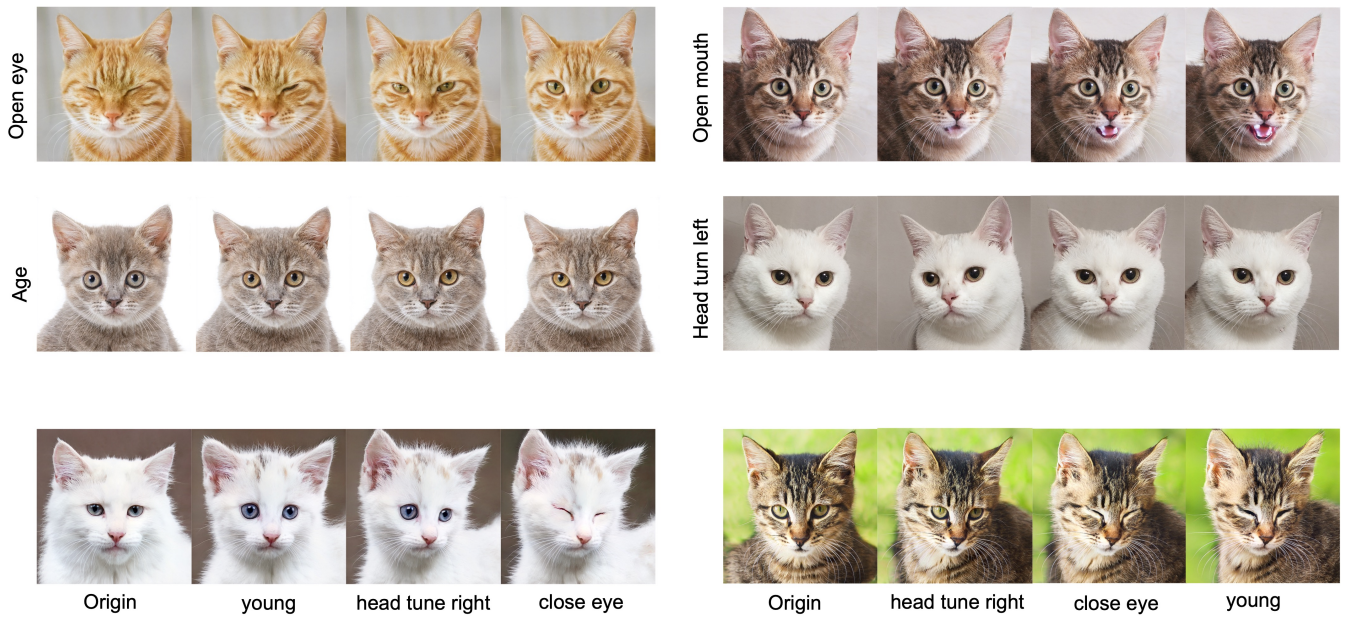


Figure 16: Results of single- and multi-attribute manipulation on animal faces.





Figure 17: Multi-attribute editing of real photographs that are reconstructed with pSp (Richardson et al. 2020) encoder.

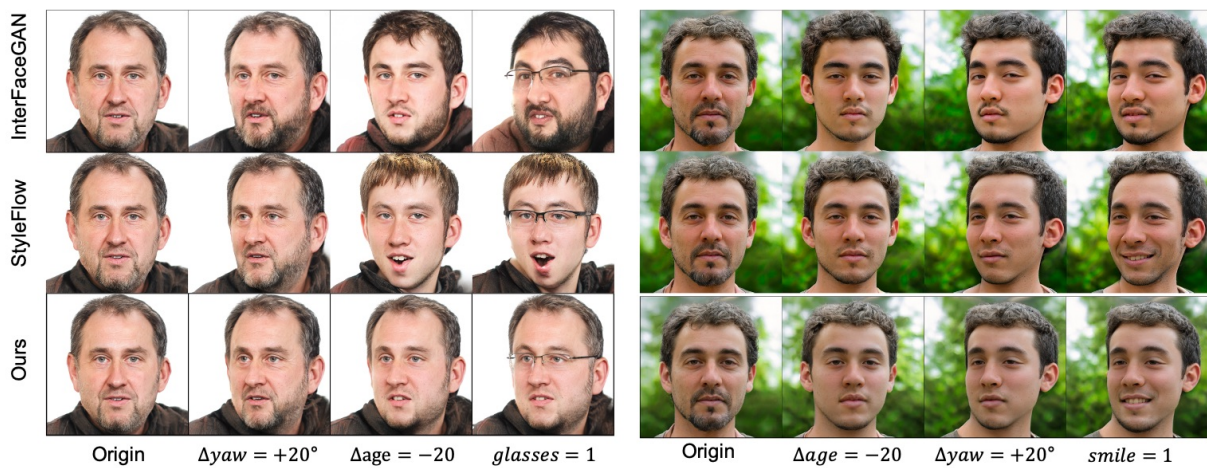


Figure 18: Comparisons of our method and competing methods in terms of multiple-attribute manipulation. As shown, the results of InterFaceGAN (Shen et al. 2020) and StyleFlow (Abdal et al. 2020) see significant identity variation when jointly manipulating multiple attribute. The identity preservation and control accuracy of yaw of InterFaceGAN and StyleFlow are problematic.



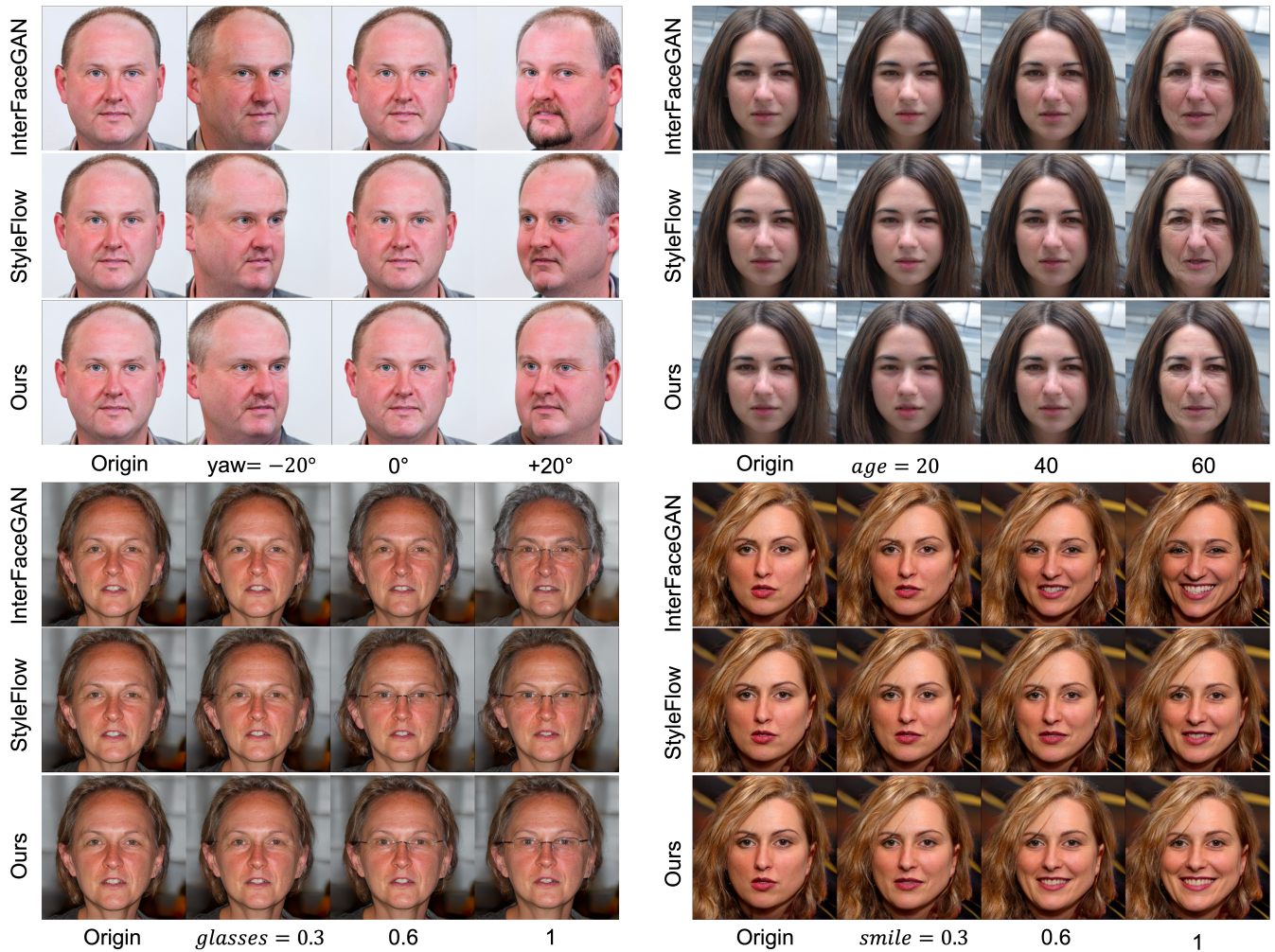


Figure 19: Comparisons of our approach and competing methods in terms of single-attribute manipulation. As shown, the InterFaceGAN (Shen et al. 2020) effectively controls the designated attribute, while the attribute disentanglement is not pleasing. E.g., the editing of yaw induces change of beards and the control of smile and glasses results in variation of age. The face identity does not hold when controlling yaw and smile. StyleFlow (Abdal et al. 2020) exhibits better attribute disentanglement. Whereas, the control precision of yaw and the smoothness of transition when adjusting the control value of smile or glasses is lower than ours.

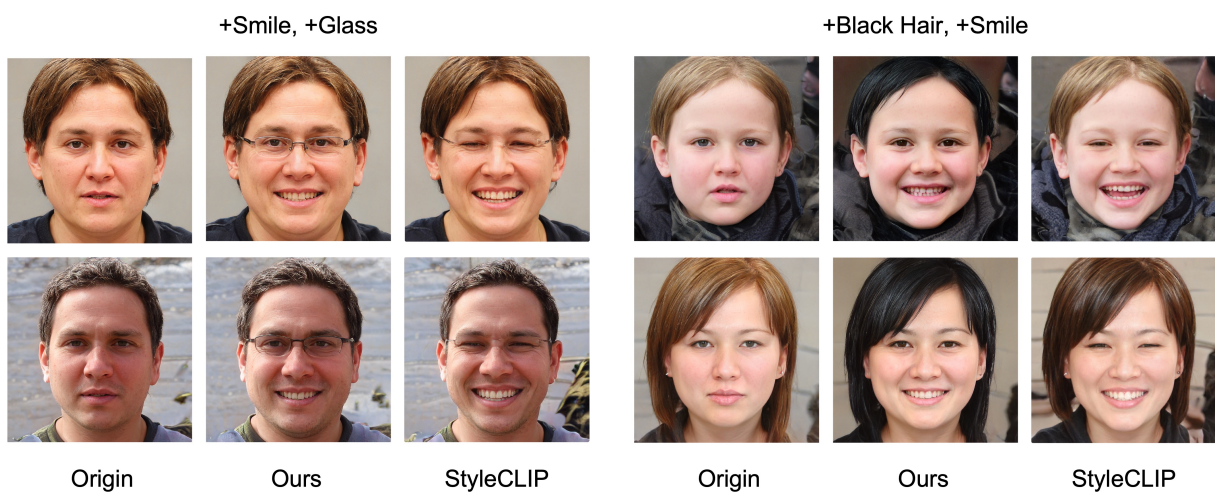


Figure 20: Comparisons of our method with StyleCLIP (Patashnik et al. 2021) in terms of multiple-attribute manipulation. As shown, the results of StyleCLIP see significant identity variation when jointly manipulating multiple attribute. The control of black hair and glasses of StyleCLIP are problematic.