# Style Transformer for Image Inversion and Editing

Xueqi Hu[1], Qiusheng Huang[1], Zhengyi Shi[1], Siyuan Li[1], Changxin Gao[3], Li Sun[1,2,*], Qingli Li[1]

[1]Shanghai Key Laboratory of Multidimensional Information Processing,
[2]Key Laboratory of Advanced Theory and Application in Statistics and Data Science,
East China Normal University, Shanghai, China
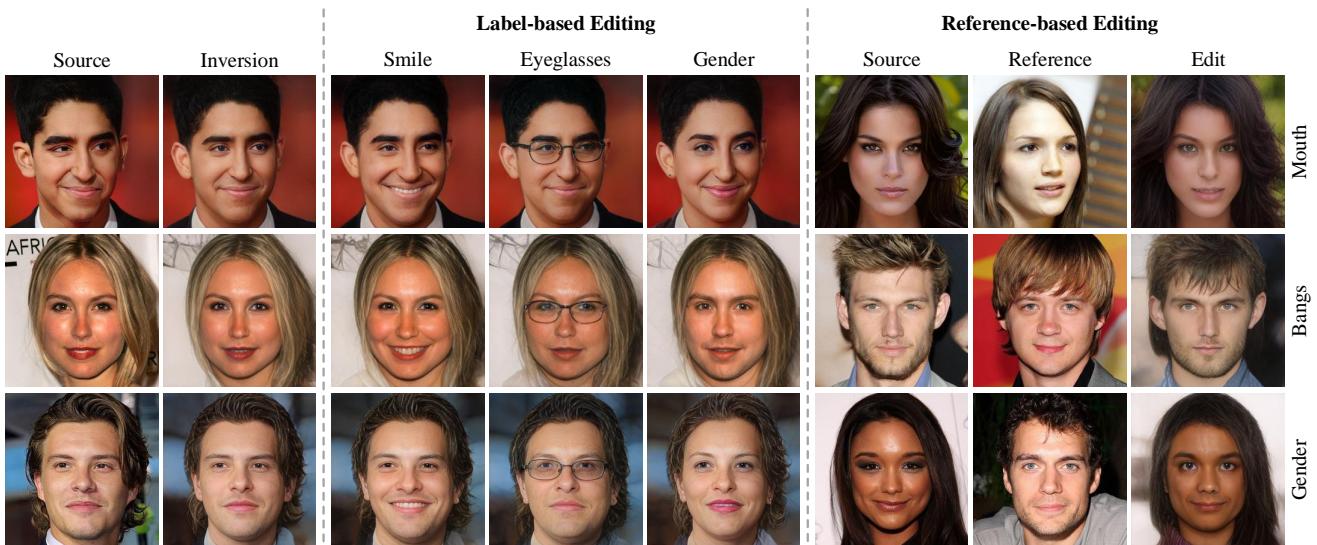[3]Huazhong University of Science and Technology, Wuhan, China

Figure 1. Image inversion and editing results of our model on CelebA-HQ dataset. From left to right, we show the inversion, label-based editing and reference-based editing. For the reference-based editing, the three columns are the source, reference and edit images, and each edit image takes the style of the reference while maintaining the source content.

## Abstract

*Existing GAN inversion methods fail to provide latent codes for reliable reconstruction and flexible editing simultaneously. This paper presents a transformer-based image inversion and editing model for pretrained StyleGAN which is not only with less distortions, but also of high quality and flexibility for editing. The proposed model employs a CNN encoder to provide multi-scale image features as keys and values. Meanwhile it regards the style code to be determined for different layers of the generator as queries. It first initializes query tokens as learnable parameters and maps them into $W^+$ space. Then the multi-stage alternate self- and cross-attention are utilized, updating queries with the purpose of inverting the input by the generator. Moreover, based on the inverted code, we investigate the reference- and label-based attribute editing through a pretrained latent classifier, and achieve flexible image-to-image translation with high quality results. Extensive experiments are carried out, showing better performances on both inversion and editing tasks within StyleGAN. Codes are available at https://github.com/sapphire497/style-transformer.*

## 1. Introduction

Generative Adversarial Network (GAN) [5, 27, 29] has been significantly improved during recent years. Particularly, with the help of AdaIN [15] or it variation ModulatedConv, StyleGAN [17, 18] is able to synthesize high resolution images with moderate quality. Therefore, utilizing

the pretrained and fixed StyleGAN for downstream tasks becomes a hot research topic, especially in the editing task of image-to-image (I2I) translation [3, 11, 31, 32, 37, 38].

To edit a given real-world image, we first need to find out its input noise vector $z$ or intermediate latent code $w$, which can faithfully reconstruct the specified real image using the pretrained generator. Then, the code is modified by an offset corresponding to the target attribute, so that it can be mapped into an edited image, while preserving the original details. Despite of the great efforts, inverting [1, 2, 30, 33, 43] or editing [3, 11, 31, 32, 37] images for StyleGAN is still challenging due to following reasons. First, there are several candidate latent embeddings. Existing methods [33, 37, 45] reveal that different choices on them are critical. Compared to $Z$ or $W$ space with a single 512-d vector, $W^+$ has the enough ability to describe image details, therefore it is suitable for inversion. In $W^+$, each image is represented by 18 different codes, and each of them is 512-d. They are given to the generator to formulate features and final synthesis from low to high resolutions in sequence. However, the code in $W^+$ can not be well edited unless imposing enough regularization. Second, the distribution in $W$ or $W^+$ are highly complex. Real images only lie on the manifold in the space [33]. Moreover, different dimensions are often entangled for a single attribute, making independent editing difficult.

This paper aims to improve the encoder-based image inversion and editing for StyleGAN at the same time. Inspired by the great success of transformer in image classification [10, 23] and object detection [6, 46], we utilize it to find the appropriate latent code in $W^+$ space for image inversion and editing tasks. The basic idea is to regard latent codes in different generation stages as query tokens, and image features at different spatial positions as keys and values, then perform the multi-head cross-attention to update the queries in an iterative way. Meanwhile, before the cross-attention, the queries are also allowed to access others through the self-attention, to enhance the regularization on them, so the final codes given to the generator become tightly linked.

Particularly, queries first interact with image features (keys) by comparing similarities between each query-key pair. Then they are organized into the attention matrix to dynamically weight the features (values) and update queries for the transformer block in next stage. The image features, used as keys and values, are obtained by a CNN encoder. To capture the image details at different resolutions, we employ a two-pyramid encoder proposed in [30] to provide multi-scale features as keys and values. Note that our model has the multiple cross-attentions from low to high resolutions, and the style queries are gradually updated by features at different scales. Therefore, general contents in queries are first formulated, and then refined by details in the higher

resolution. After several times self- and cross-attentions, queries absorb enough details from the input image, so they can be utilized to invert it by the pretrained generator.

We are further interested in the way to edit the codes for translating a specified attribute. Traditional approaches [11, 31, 37] assume the linear separations in the latent space for a binary attribute, so inverted code from different images are edited by the same direction. We argue that the identical direction is not optimal for the editing quality, and may reduce the result diversity. Inspired by [7, 14], we divide the image editing in StyleGAN into two different types. One is label-based editing, in which only target label is specified. The other is reference-based editing, which requires another image to supply the desired style. For the former, a pretrained non-linear latent classifier is used to determine the direction. it computes a loss for the inverted code according to the target label, and its gradient is back-propagated to the code, giving the editing direction. In the latter case, we want to determine the exact editing vector from the reference. Therefore, the inverted code from the source is used as query, and from the reference as key and value. The cross-attention is performed between them. The parameters in the attention module is trained under the supervision from the latent classifier, encouraging the edited attribute to be similar with the reference and other attributes without any changes. The proposed editing method is able to give the diverse results while maintaining the quality of image.

The contribution of the paper is summarized into following aspects. First, we propose novel multi-stage style transformer in $W^+$ space to invert image accurately. The transformer includes the self- and cross-attention modules, in which the style queries gradually get updated from the multi-scale image features. Second, we characterize the image editing in StyleGAN into label-based and reference-based, and use a non-linear classifier to generate the editing vector. Diverse and fidelity editing results are obtained.

## 2. Related Works

**GAN inversion** is first proposed in [44] and becomes important due to the wide applications of some recent generators. There are basically two ways, either encoder-free or encoder-based. The former does not have any training parameters, and the latent codes are directly optimized by the gradient mainly from reconstruction loss. To deal with the complex latent structure, Abdal *et al*. [1] invert a real image in $W^+$ space, and use pixel-wise MSE and perceptual loss with Adam [20] to tune the code. Image2StyleGAN++ [2] extends the code space to the layerwise additive noise vector to decrease the distortion. Although such a method can reliably find the code through multi-step iterations, it is inefficient and its code lacks the editability.

In contrast, the encoder-based method intends to train a common model to achieve inversion for all images. It

improves the editing ability and is efficient during inference. IDInvert [43] utilizes a CNN as encoder to output the code. Except the reconstruction and perceptual loss, it is trained by an extra adversarial loss. pSp [30] designs a two-pyramid encoder to provide multi-scale features, and maps them to the style vector through multiple convolution layers. Benefiting from strong features, pSp achieves less distortion. ReStyle [4] uses the encoder to give the residual style to refine the inversion in the iterative way. E4e [33] analyzes the distortion-editability trade off for inversion and editing tasks in $W^+$. It sacrifices the inversion accuracy to improve the editablity, constraining the codes for different layers close with each other. Kim *et al.* [19] and Wang *et al.* [36] depart from $W^+$ space and enhance the code with spatial dimensions, so that more information are given to the generator to lower down the distortions. Compared to previous works, our method lies strictly in $W^+$, and it is able to achieve minimal distortion and high quality editing at the same time.

**Latent code manipulation** for pretrained StyleGAN is often used to edit the attribute and achieve I2I translation, either in the supervised or unsupervised manner. GANSpace [11] and Sefa [32] adopt PCA to find the principal directions in $W$ space. They are responsible for controlling the pose, gender or background. Note that for a particular attribute, these works specify the same direction on all latent codes to realize editing. Voynov and Babenko [34] train a simple module to edit the input, and use a reconstructor in pixel domain to interpret the editing, finding noticeable directions. LatentCLR [40] builds a learnable direction model to edit the code and uses contrastive loss to train it. Hence, these two models give different images unique editing directions. All the above works are unsupervised, without requiring attribute label for editing. But only limited directions for some attributes can be found.

The supervised methods can identify directions for more attributes, especially the local ones. InterfaceGAN [31] trains the linear binary SVM in latent space to obtain a separation plane, whose normal vector controls its corresponding attribute. StyleSpace [37] finds the control direction in a precise way guided by the semantic mask. Moreover, they propose to edit the code in $S$ space which is defined by the affine layer after $W$. Wang *et al.* [35] further extend $S$ space by tracing back the gradient flow to its previous stage, making the change more accurately. Note that these works still share the same editing direction for all images. Recently, StyleFlow [3] conditionally manipulates the images using Continuous Normalizing Flows (CNF). Yao *et al.* [38] propose a latent transformation module to generate adaptive directions for different images. Wang *et al.* [36] utilize a CNN encoder to provide multi-scale features to supplement the $1 \times 1$ style vector, which actually adapts directions at different locations.

However, previous works only deal with the label-based editing. Collins *et al.* [8] apply k-means clustering on features to obtain channel-wise masks, determining which channels are locally semantic-aware. The cluster memberships of the reference further guide local attribute editing for the source. Different from above works, our work is strictly in $W^+$, and we realize both the label- and the reference-based editing.

## 3. Framework of Style Transformer

We aim to achieve accurate image inversion for Style-GAN by our proposed style transformer in $W^+$ space. Given a real image $I \in \mathbb{R}^{H \times W \times 3}$, our model is able to specify $N$ different style vectors denoted by $w_n \in \mathbb{R}^{512}$, where $n = 1, 2, \cdots, N$ is the index of the vector injected into the different stages of the generator $G$. For simplicity, we use $w \in \mathbb{R}^{N \times 512}$ without any index to represent all $w_n$. Note that in StyleGAN2, $w_n$ is first projected by the affine layer $A$, then it affects the corresponding layers by modulating on the convolution kernels.

Fig. 2 illustrates the overview of the proposed framework. The input image $I$ is encoded by $E$, generating a series of image features $F_1$ to $F_3$ in multi-resolutions [30]. $N$ different queries, output from an MLP, access these features through the transformer blocks in a sequential way, forming the final code $w$ for the generator. The initial input $z_n \in \mathbb{R}^{512}$ to the MLP is also learnable, and it is gradually updated into $w$, which is suitable for inverting $I$. By training all parameters including the transformer blocks, the encoder $E$, the MLP and the initial $z_n$, the pretrained $G$ can utilize the final $w$ to reconstruct input $I$ with minimal distortions.

### 3.1. Style Transformer Block

The style transformer block is the key component for image inversion. The same structure is applied for 3 times in the model, exploiting the image details from $F_1$ to $F_3$, respectively. The specific design within the block is shown on the right of Fig. 2. Basically, there are two types of attention, which are multi-head self-attention and cross-attention. In addition, we follow the design routines for transformer, incorporating the residual connection, the normalization and FFN module into the block.

**Style query initialization.** Given a single style code $w$, high fidelity image can be synthesized by StyleGAN generator. However, $W^+$ space needs $N$ different style vectors to reconstruct one image, and they essentially describe the details at different scales, therefore are employed to affect features of different resolutions in the generator. A common choice in transformer decoder is to randomly initialize beginning query tokens and keep them as learnable parameters in the model. However, considering the fact that code distribution in $W$ space is complex and far from Gaussian
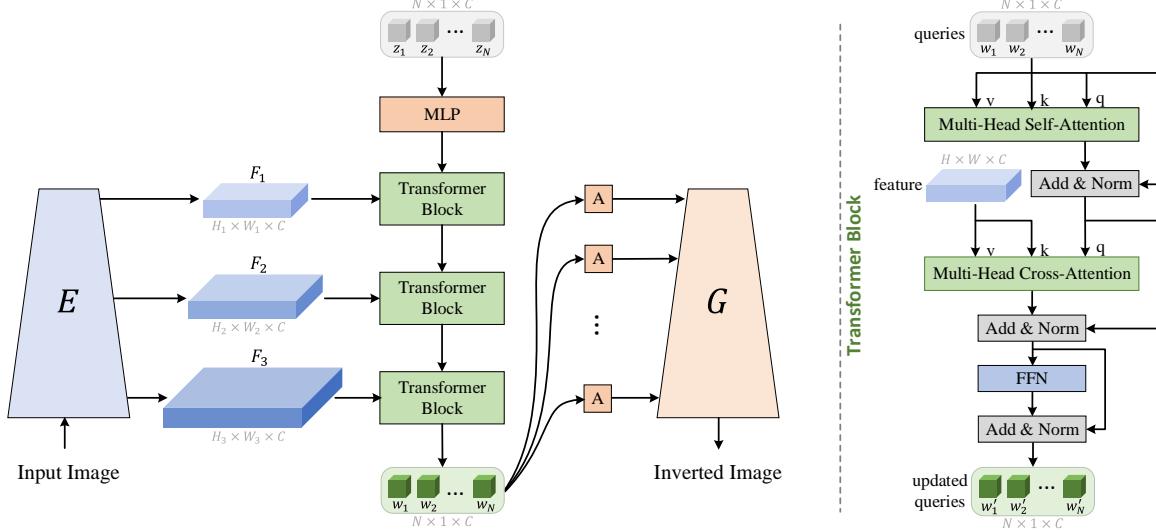
3

Figure 2. The overall framework of Style Transformer for image inversion. We build the multi-stage transformer-based model to update the code in $W^+$ space. Details within the transformer block are depicted on the right. Each has a multi-head self- and cross-attention block, following the common routine in transformer model.

prior, we utilize the pretrained MLP in StyleGAN to first map each individual code $z_n$ to the beginning style query $w_n$ in $W$ space, and then update $w_n$ through the self and cross-attention operations. Note that $z_n \sim N(0, I)$ is sampled from standard Gaussian, and set as model parameters. Moreover, the pretrained MLP is finetuned during training.

**Multi-Head Self-Attention.** The self-attention is performed among $N$ different query tokens $q_1, q_2, \cdots, q_N$. It intends to find the potential relation between any pair of them, and route the value to connect them. We denote all of q as $X_q \in \mathbb{R}^{N \times 512}$. The query $Q$, key $K$ and value $V$ are all projected from $X_q$ according to Eq. (1). Note that $W_Q^{self}$, $W_K^{self}$ and $W_V^{self} \in \mathbb{R}^{512 \times 512}$ are learnable projection heads in the self-attention module, which do not change the feature dimension.

$$Q = X_q W_Q^{self}, \quad K = X_q W_K^{self}, \quad V = X_q W_V^{self} \quad (1)$$

The multi-head attention operation is formulated as in Eq. (2), where $Q_i$, $K_i$ and $V_i$ are query, key and value in the $i$th head, and $Attn$ is result from that head. The feature dimension $d = 512/H$, and $H$ is the number of attention heads.

$$Attn(Q, K, V) = \text{Softmax}(\frac{Q_i K_i^T}{\sqrt{d}})V_i \quad (2)$$

The final update on $X_q$ from the self-attention is $MHA$ in Eq. (3). $W^o \in \mathbb{R}^{512 \times 512}$ is also learnable, being responsible for fusion the results $Attn$ from different heads.

$$MHA(Q, K, V) = [Attn(Q_i, K_i, V_i)]_{h=1:H} W^o \quad (3)$$

**Multi-Head Cross-Attention.** The self-attention does not involve any image feature in its computation. Therefore,

we further design the cross-attention for the inversion task, so that the query tokens can obtain information from image features $F_1$, $F_2$ and $F_3$ in different resolutions. In the cross-attention, features of key and value are from the encoder $E$, while queries are computed by the linear projection on the previous results from self-attention block. Particularly, we have the query, key and value according to Eq. (4), where $W_Q^{crs}$, $W_K^{crs}$ and $W_V^{crs}$ share the similar settings with the self-attention.

$$Q = X_q W_Q^{crs}, \quad K = F_i W_K^{crs}, \quad V = F_i W_V^{crs} \quad (4)$$

The multi-head cross-attention is carried out in the same way as is shown in Eq. (2) and Eq. (3). After that, the updated query tokens are given to an FFN to refine itself, and the results are further passed to the transformer block in the next stage, mining the details from finer resolution features.

### 3.2. Training Objectives for Image Inversion

During training, the backbone $G$ of StyleGAN (including the affine layer $A$) is strictly fixed. The gradients from the loss only tune other parameters. Note that we use the same training objectives as pSp [30]. Particularly, to give the accurate reconstruction, the $L_2$ loss between the input $I$ and its inverted version $\hat{I}$ from $G$ is calculated. Meanwhile, LPIPS [42], a similarity metric, which is computed based on the features in an Inception net $F(\cdot)$, is also adopted, specifying another objective $L_{LPIPS} = \|F(I) - F(\hat{I})\|_2$. Additionally, to keep the identity of the inverted image, we incorporate a pretrained ArcFace model [9] $R(\cdot)$ for the ID loss $L_{ID} = 1 - \langle R(I), R(\hat{I}) \rangle$, so that the cosine similarity of $I$ and $\hat{I}$ can be maximized. Notice that we do not adopt any adversarial loss during training.
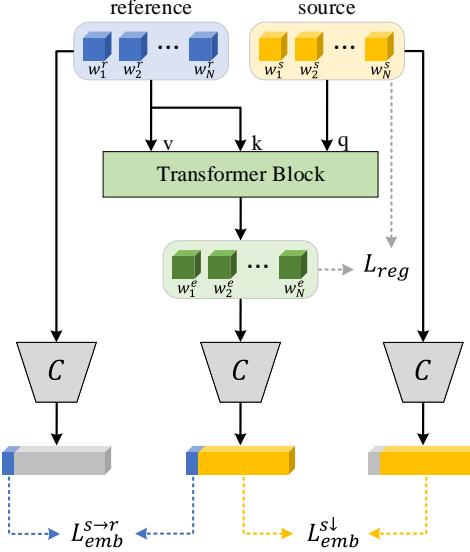
Figure 3. Reference-based editing module and its training strategy. The inverted codes $w^s$ and $w^r$, from the source and reference, are given to the transformer module $T$, specifying the code $w^e$ for edited image. $C$ is an attribute classifier in $W^+$, by which we constrain the editing attribute being similar with the reference, while others staying the same as the source.

## 4. Image Editing in Style Transformer

Image editing by the fixed StyleGAN is an important application not only for itself, but also for evaluating the quality of image inversion. The low distortion is only the one aspect, the flexible and high fidelity editing are also important. As is described in [7, 14], there are two types of editing, either through the target label or a reference image in the desired domain. Previous works [31, 37, 38] focus on the former, but few works deal with the referenced-based editing, which potentially provides diverse results.

Typically, given the inverted style code $w^s \in \mathbb{R}^{N \times 512}$ for the source $I_s$, and a desired target attribute, we need to determine an offset $\Delta w$, so that $w^e = w^s + \Delta w$ can be mapped to an edited image $\tilde{I} = G(w^e)$ with the desired attribute different from $I_s$, but keep content of $I_s$. In the reference-based editing, another image $I_r$ is given as an extra input. Since our style transformer can invert image with negligible distortion, we train a latent classifier $C$ for $K$ binary attributes in $W^+$ space to guide the editing like [38]. Concretely, given a code $w$ inverted from an image, the classifier computes several embedding features $C_f^k$ corresponding to the $k$th attributes, and the final logits $C_l^k$ for the BCE loss $L_{bce}$. During editing, $C$ is fixed to evaluate $w^e$.

### 4.1. Reference-based Editing

**Module design.** We design a simple module $T$ to translate a particular attribute according to the inverted code $w^r$

from reference $I_r$. Since both $w^s$ and $w^r$ represent images with almost no distortions, these codes contain enough information about the edited attribute. $T$ should be able to specify $\Delta w$ based on $w^r$ and $w^s$. Again, a cross-attention structure is chosen, as is shown in Fig. 3. $Q = w^s W_Q^{edt}$ are used as a series of query tokens, while $K = w^r W_K^{edt}$ and $V = w^r W_V^{edt}$ are key and value tokens, projected from $w^r$. According to [31,37], some local attributes are only depended on a single $w_n$ for a particular resolution in $G$. So we choose a routing scheme in [24] different from Eq. (2). The idea is to make $\text{Softmax}_Q$ normalize over queries not keys. Then re-norm the matrix over keys by $\text{Norm}_K$ as is Eq. (5). This strategy assigns the value feature $V$ to queries in the unique way, so that a value token from $w^r$ affects only a few tokens in $w^s$.

$$T(w^s, w^r) = \Delta w = \text{Norm}_K(\text{Softmax}_Q(\frac{QK^T}{\sqrt{d}}))V \quad (5)$$

**Loss designs.** To guarantee the $k$th attribute editing results, we design the following loss terms to train the projection head in $T$. Particularly, we constrain the code $w^e$ after editing by $L_{emb}^{s \to r}$ as is shown in Eq. (6):

$$L_{emb}^{s \to r} = \|C_f^k(w^e) - C_f^k(w^r)\|_2 \quad (6)$$

Here $C_f^k$ is the $k$th attribute embedding from the pretrained latent classifier $C$. $L_{emb}^{s \to r}$ ensures the edited attribute to be similar with $I^r$. At the same time, other attributes denoted by $\cancel{k}$ should stay close with the source $I_s$, giving $L_{emb}^{s\downarrow}$ in Eq. (7):

$$L_{emb}^{s\downarrow} = \|C_f^{\cancel{k}}(w^e) - C_f^{\cancel{k}}(w^s)\|_2 \quad (7)$$

Finally, we regularize $L_{reg} = \|\Delta w\|_2 = \|w^e - w^s\|_2$, so edited image $\tilde{I}$ does not change much.

### 4.2. Label-based Editing

Compared to reference-based, label-based editing is relatively easy. So we adopt an encoder-free method to edit $w$ based on the latent classifier $C$. We emphasize that for each $I_s$, there should be a unique direction $n_{\Delta w}^k$ for the $k$th attribute editing, which is determined by the gradient back-propagated from the classifier $C$. Note that the first-order gradient on $w$ is $g = \nabla_w L_{bce}(C_l^k(w^s), y_t)$, and the direction becomes $n_{\Delta w}^k = -g/\|g\|_2$. Here $y_t$ is the target label, and $C_l^k(w^s)$ is the logits after sigmoid.

We also investigate the method based on the second-order derivative $H$, which is the Hessian matrix. Similar with [26], a randomly sampled unit vector $d$ is first obtained, then it is scaled by a small number $\xi$. Then we evaluate the Hessian vector product by Eq. (8). According to power iteration, $d \leftarrow Hd$ converges to the dominant eigenvector, so we let $g = Hd$.

$$Hd \approx \frac{\nabla_r L_{bce}(C^k(w^s+r), y_t)|_{r=\xi d} - \nabla_r L_{bce}(C^k(w^s), y_t)|_{r=0}}{\xi} \quad (8)$$
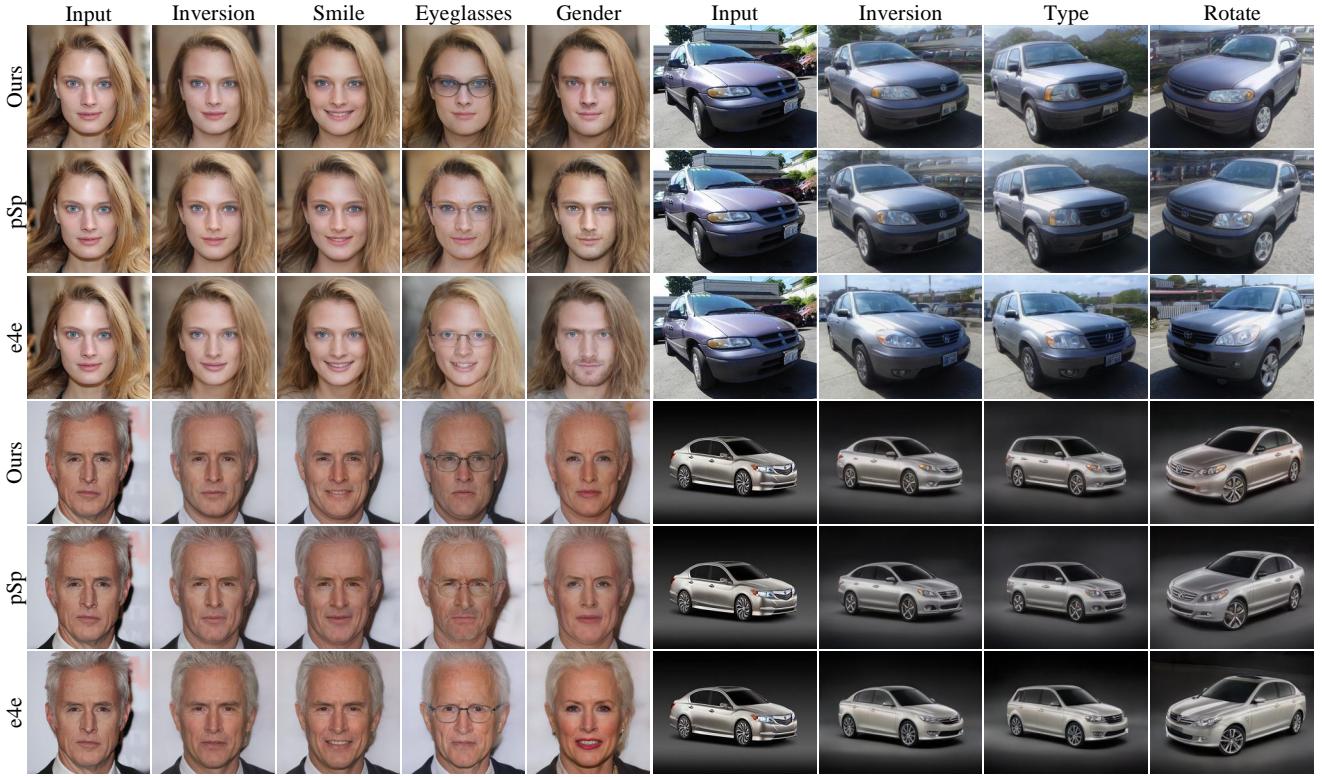
Figure 4. Qualitative results of image inversion. Our method is compared with pSp and e4e. Besides the inverted image, we also list images edited by InterFaceGAN [31] for faces. For cars, we use the directions provided in GANSpace [11] for editing.

| Domain | Method | Inversion | | | | Editing | | Model Size | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | MSE↓ | LPIPS↓ | FID↓ | SWD↓ | FID↓ | SWD↓ | Params(M)↓ | FLOPs(G)↓ | Time(s)↓ |
| Face | pSp | 0.037 | 0.169 | 31.52 | 15.07 | 46.64 | 29.05 | 267.3 | 72.55 | 0.0668 |
| | e4e | 0.050 | 0.209 | 36.16 | 17.25 | 47.45 | 25.10 | 267.3 | 72.55 | 0.0659 |
| | Ours | **0.036** | **0.166** | **28.31** | **14.00** | **40.57** | **23.21** | **40.6** | **36.37** | **0.0436** |
| Car | pSp | 0.115 | 0.298 | 17.24 | 19.76 | 27.25 | 36.01 | 238.0 | 66.11 | 0.0565 |
| | e4e | 0.110 | 0.314 | 14.68 | 18.25 | 21.50 | 27.57 | 238.0 | 66.11 | 0.0541 |
| | Ours | **0.089** | **0.245** | **13.58** | **16.14** | **21.24** | **25.28** | **40.6** | **36.34** | **0.0435** |

Table 1. Quantitative comparison for different inversion methods. To consider the distortion-editability trade-off, we list metrics for image editing to give a comprehensive evaluation on them. We also list the parameters and FLOPs of the three methods, *Time* means the inference time of an iteration.

## 5. Experiments

### 5.1. Implementation Details

All experiments are implemented on StyleGAN2 [18] pretrained on FFHQ [17] and LSUN Cars [39] datasets. We build our model based on the pSp encoder for multi-scale image feature. For face domain, we train the inversion model on FFHQ dataset and evaluate on CelebA-HQ [16] test set. For car domain, the inversion model is trained and evaluated on Stanford Cars [21] dataset. The synthesis network in StyleGAN2 is fixed and all other parameters in our model is trainable.

### 5.2. Inversion Results

We compare our model with pSp [30] and e4e [33], which are two state-of-the-art encoder-based inversion methods. Qualitative and quantitative results are shown in Fig. 4 and Tab. 1. Our model is validated in three aspects: the perceptual quality of inversion, the ability of editing, and the model size. MSE and LPIPS evaluate the pixel and perceptual similarity of input and inverted images. FID [13] and SWD [28] measure the distance between two distributions of real and generated images, indicating the visual
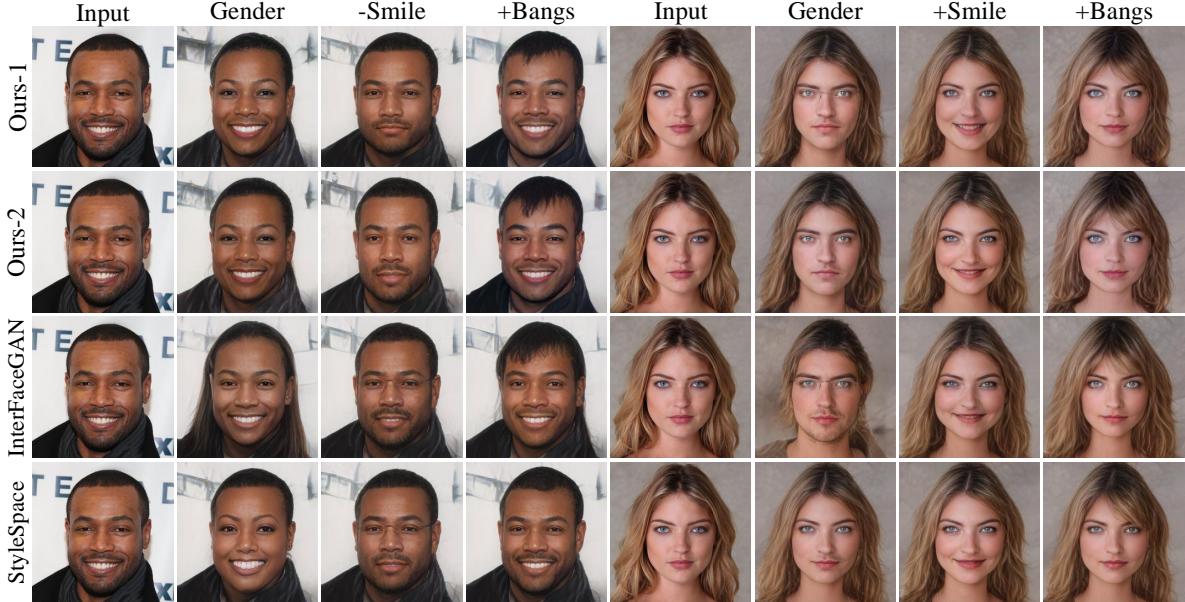
Figure 5. Qualitative comparison on different label-based editing methods. We list the results for editing "Gender", "Smile" and "Bangs", and compare them with [31, 37]. Note that we evaluate both the first- and second-order editing methods proposed in Sec. 4.2.

| Method | Gender | | Smile | | Bangs | |
|---|---|---|---|---|---|---|
| | FID↓ | SWD↓ | FID↓ | SWD↓ | FID↓ | SWD↓ |
| InterFaceGAN | 48.72 | 19.43 | 40.03 | 18.94 | 44.01 | 29.41 |
| StyleSpace | 37.31 | 17.31 | 34.72 | 15.46 | 42.91 | 20.96 |
| Ours-1 | 38.73 | 17.83 | 33.50 | **14.89** | 41.15 | 19.30 |
| Ours-2 | **34.84** | **16.14** | **32.88** | 15.23 | **40.14** | **18.53** |

Table 2. Quantitative comparison of label-based editing on three attributes, corresponding to Fig. 5.
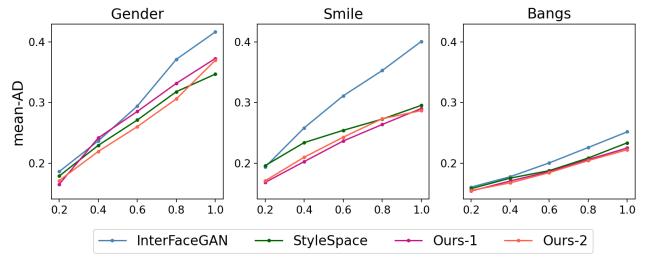


Figure 6. Mean-AD results of label-based editing on three attributes compared with [31, 37], lower means better. Ours-1 and Ours-2 represent our first- and second-order methods, respectively.

quality of generated images. To compare the editing ability of three methods, we adopt InterFaceGAN [31] in face domain to edit the latent codes generated by each method. For car domain, we apply GANSpace [11] to find the semantic directions. The metrics are averaged over the editing results of the whole test set. Our model is outperformed in inversion and of higher editability. Moreover, we list the parameters, FLOPs and inference time of three methods in Tab. 1. Compared with Convnet, the transformer used in our model has only 18 or 16 tokens for face and car domain, hence it is lightweight and efficient.

## 5.3. Editing Results

We apply reference- and label-based editing on CelebA-HQ dataset, in which each image has the label of 40 facial attributes. We invert images to latent codes using our pretrained inversion model, and train a 40-class latent classifier. The latent classifier consists of 4 fully-connected layers, in which there is an independent branch for each attribute before the prediction, leading to the independent embedding feature.

**Label-based Editing.** We first apply our pretrained inversion model to obtain the latent codes of images, and use the first- and second-order methods to edit the images to have the target attributes. Desirable results can be generated by only ONE iteration. We evaluate first- and second-order methods illustrated in Sec. 4.2 and compare our results with InterFaceGAN [31] and StyleSpace [37]. Qualitative results and metrics are shown in Fig. 5 and Tab. 2. Moreover, we measure the disentanglement of attributes by calculating the Attribute Dependency (AD) [37], which indicates the degree of changes in other attributes while editing one attribute. A multi-branch attribute classifier based on ResNet-50 [12] is applied to obtain the predicted logits of images. We measure the changes $\Delta l$ between the input and edited images, and normalize $\Delta l$ by $\sigma(l)$, which is the standard deviation computed from the logits of numerous generated images. For a target attribute $k$, we calculate the mean-AD
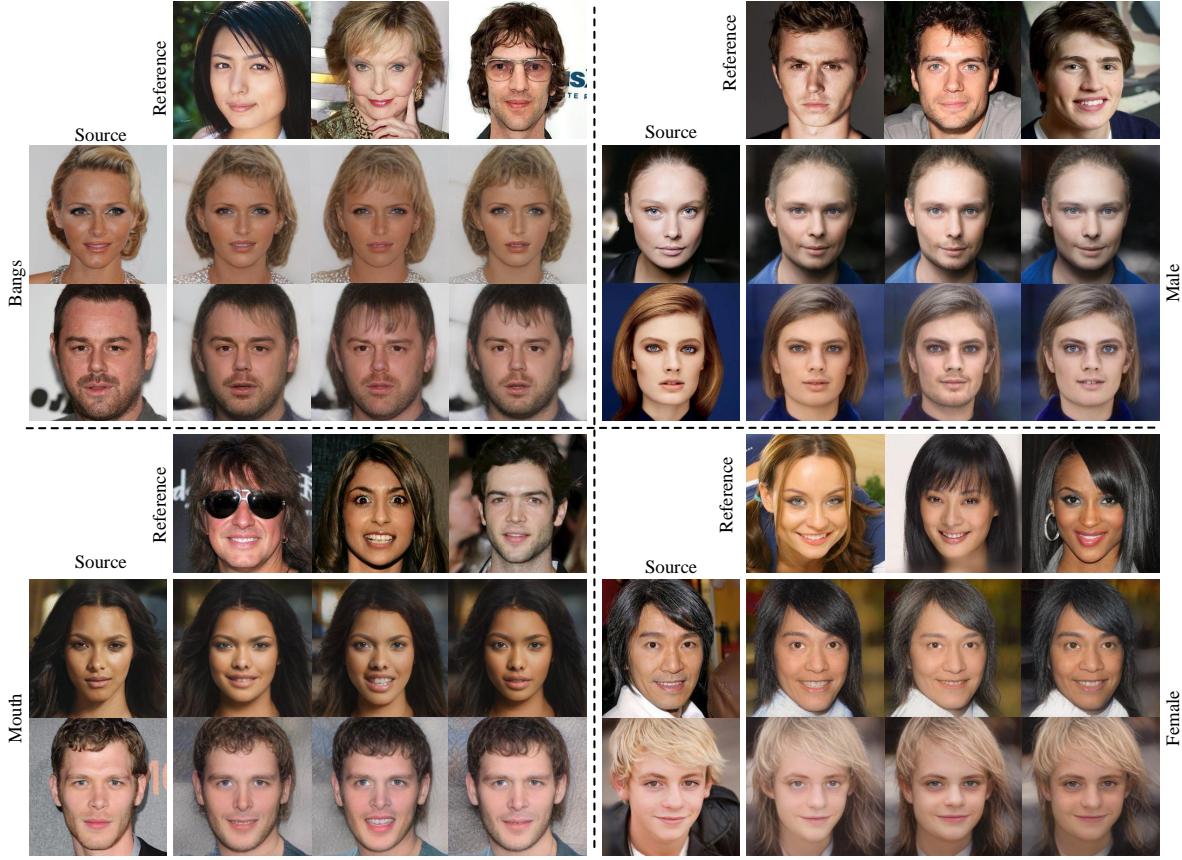
Figure 7. Reference-based editing results. Given a pair of source and reference images, we first utilize the proposed method to find their inverted codes in $W^+$. Then the transformer block described in Sec. 4.1 is used to take the "bangs", "mouth" and "gender" style in the reference code, and apply them onto the source.

| Method | Quality(%) | | | Disentanglement(%) | | |
|---|---|---|---|---|---|---|
| | **BA** | **GE** | **GO** | **BA** | **GE** | **GO** |
| InterFaceGAN | 15.00 | 7.50 | 9.17 | 11.67 | 1.67 | 8.33 |
| StyleSpace | 10.83 | 10.00 | 13.33 | 18.33 | 15.00 | 10.83 |
| Ours-1 | 25.83 | 39.17 | 31.67 | **35.83** | 34.17 | 30.00 |
| Ours-2 | **48.33** | **43.33** | **47.50** | 34.17 | **49.17** | **49.17** |

Table 3. User study of label-based editing compared with [31], [37]. **BA**, **GE** and **GO** represent 'Bangs', 'Gender' and 'Goatee' attributes.

on other attributes $\not{k}$ as $\mathbb{E}(\Delta l_i / \sigma(l_i))$, where $i \in \not{k}$ is the index of fixed attributes. Fig. 6 illustrates the mean-AD over the degree of target attribute changes $\Delta l_k / \sigma(l_k)$, Our methods perform better than InterFaceGAN and StyleSpace, and the second-order method is of higher disentanglement.

Considering human judgements, we further conduct a user study. We ask 60 volunteers to evaluate the methods in two aspects: image quality and disentanglement. Results are shown in Tab. 3. The detailed algorithms of first- and second-order methods are provided in Appendix C.

**Reference-based Editing.** The reference-based editing module is trained for different attributes individually. To ensure the module takes the style from reference image, we randomly divide the training images into source and reference sets, instead of depending on the labels. We train the module on three attributes and qualitative results are shown in Fig. 7. The edited images take the relevant attributes from different reference images, and they appear the similar style on the translated attribute. Note that the reference-based editing module is trained only in the latent space, resulting in less diversity compared with directly editing on images. Whereas, different from the optimized-based method [8], our model can commonly apply to all images, which is lightweight and more flexible.

## 5.4. Ablations and Analysis

We further validate the benefit of transformer by comparing among pSp [30], our full model with both self- and cross-attention and ours w/o self-attention in Tab. 4. [30] maps image features to $w+$ by individual mapping networks, though $w+$ obtain the image features directly and completely, the relation between each $w$ is not tightly

| Method | MSE↓ | LPIPS↓ | Params(M)↓ | FLOPs(G)↓ | Time(s)↓ |
|---|---|---|---|---|---|
| pSp | 0.0373 | 0.1693 | 267.3 | 72.55 | 0.0668 |
| Ours w/o self | 0.0369 | 0.1716 | **37.3** | **36.31** | **0.0429** |
| Ours full | **0.0363** | **0.1665** | 40.6 | 36.37 | 0.0436 |

Table 4. Ablations of transformer structure. *Time* means the inference time of an iteration. The best results are indicated in **Bold**.

enough. In our model, cross-attention is necessary to update queries by fusing image features, and self-attention is also important in constructing the potential relation between queries.

## 6. Limitations

We now discuss limitations, which we have already realized, for our work. First, for the inversion task, although our proposed method achieves improved reconstruction quality, there are still some differences between the input and reconstructed images, especially for the out-of-domain input. We think it is mainly caused by the finite discriminative ability of $W^+$ space. As is described in [36], the distortion can be significantly reduced by adding more information from the source. Moreover, since we apply the multi-head attention, the training speed is slower due to the complex matrix multiplication. Second, for the reference-based editing task, we adopt a transformer-based module in the latent space, resulting in less diversity for some attributes compared with direct editing on the images, in which the mode seeking loss [25] can encourage the diversity in the pixel domain. But our method is lightweight and more flexible.

## 7. Conclusion

This paper presents a transformer-based image inversion and editing method for StyleGAN. We choose $W^+$ space to represent real images, which needs to determine multiple style codes for different layers of the generator. To effectively exploit information from input image, we design a multi-stage transformer module, which mainly consists of the self- and cross-attention. In the initial stage, the MLP maps a set of learnable noise vectors into the codes in $W^+$, and then they are iteratively updated by the two types of attention operations, so the codes from the final stage can reconstruct the input accurately. Based on them, we are able to carry out label- and reference-based editing in a flexible way. Given a required label, an encoder-free strategy is employed to find the unique editing vector according to the gradient from a pretrained latent classifier. Meanwhile, given a reference code, a transformer block is trained to edit the source, so that the result takes the relevant style from the reference. Experiments show the proposed image inversion and editing method achieves less distortions and higher quality at the same time.

## References

[1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, 2019. 2

[2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8296–8305, 2020. 2

[3] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *ACM Transactions on Graphics (TOG)*, 40(3):1–21, 2021. 2, 3

[4] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder via iterative refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6711–6720, 2021. 3

[5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 1

[6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020. 2

[7] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8188–8197, 2020. 2, 5

[8] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5771–5780, 2020. 3, 8

[9] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019. 4

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2

[11] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. *arXiv preprint arXiv:2004.02546*, 2020. 2, 3, 6, 7

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 7

[13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a

two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6

[14] Qiusheng Huang, Zhilin Zheng, Xueqi Hu, Li Sun, and Qingli Li. Bridging the gap between label-and reference-based synthesis in multi-attribute image-to-image translation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14628–14637, 2021. 2, 5

[15] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017. 1

[16] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 6

[17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1, 6

[18] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 1, 6, 12

[19] Hyunsu Kim, Yunjey Choi, Junho Kim, Sungjoo Yoo, and Youngjung Uh. Exploiting spatial dimensions of latent in gan for real-time image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 852–861, 2021. 3

[20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 2, 12

[21] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013. 6

[22] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019. 12

[23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*, 2021. 2

[24] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *arXiv preprint arXiv:2006.15055*, 2020. 5

[25] Qi Mao, Hsin-Ying Lee, Hung-Yu Tseng, Siwei Ma, and Ming-Hsuan Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1429–1437, 2019. 9

[26] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018. 5

[27] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*, 2017. 1

[28] Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot. Wasserstein barycenter and its application to texture mixing. In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 435–446. Springer, 2011. 6

[29] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In *Advances in neural information processing systems*, pages 14866–14876, 2019. 1

[30] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2287–2296, 2021. 2, 3, 4, 6, 8, 13

[31] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 2, 3, 5, 6, 7, 8, 12

[32] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1532–1540, 2021. 2, 3

[33] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–14, 2021. 2, 3, 6, 13

[34] Andrey Voynov and Artem Babenko. Unsupervised discovery of interpretable directions in the gan latent space. In *International Conference on Machine Learning*, pages 9786–9796. PMLR, 2020. 3

[35] Rui Wang, Jian Chen, Gang Yu, Li Sun, Changqian Yu, Changxin Gao, and Nong Sang. Attribute-specific control units in stylegan for fine-grained image manipulation. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 926–934, 2021. 3

[36] Tengfei Wang, Yong Zhang, Yanbo Fan, Jue Wang, and Qifeng Chen. High-fidelity gan inversion for image attribute editing. *arxiv:2109.06590*, 2021. 3, 9

[37] Zongze Wu, Dani Lischinski, and Eli Shechtman. Stylespace analysis: Disentangled controls for stylegan image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12863–12872, 2021. 2, 3, 5, 7, 8

[38] Xu Yao, Alasdair Newson, Yann Gousseau, and Pierre Hellier. A latent transformer for disentangled face editing in images and videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13789–13798, 2021. 2, 3, 5

[39] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a

large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 6

[40] Oğuz Kaan Yüksel, Enis Simsar, Ezgi Gülperi Er, and Pinar Yanardag. Latentclr: A contrastive learning approach for unsupervised discovery of interpretable directions. *arXiv preprint arXiv:2104.00820*, 2021. 3

[41] Michael R Zhang, James Lucas, Geoffrey Hinton, and Jimmy Ba. Lookahead optimizer: k steps forward, 1 step back. *arXiv preprint arXiv:1907.08610*, 2019. 12

[42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 4

[43] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image editing. In *European conference on computer vision*, pages 592–608. Springer, 2020. 2, 3

[44] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *European conference on computer vision*, pages 597–613. Springer, 2016. 2

[45] Peihao Zhu, Rameen Abdal, Yipeng Qin, John Femiani, and Peter Wonka. Improved stylegan embedding: Where are the good latents? *arXiv preprint arXiv:2012.09036*, 2020. 2

[46] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 2

# Appendix

## A. Training Details

We adopt a pretrained StyleGAN2 [18] generator in our experiments, in which the synthesis network is fixed and the mapping network (MLP) is trained. In the multi-head attention of the transformer block, the number of heads is set to 4, and the dimension of each head is 512. For inversion task, the Ranger optimizer is used in training, which is a combination of Rectified Adam [22] with the Lookahead technique [41]. We train the model for $6 \times 10^5$ iterations with a batch size of 8, the learning rate is set to $1 \times 10^{-4}$. For the reference-based editing task, we use the Adam [20] optimizer to train the model for $1 \times 10^4$ iterations with a batch size of 8, the learning rate is set to $1 \times 10^{-3}$. All experiments are implemented on 2 NVIDIA RTX 2080Ti GPUs.

## B. Label-based Editing Methods

We propose first- and second-order label-based editing methods in the main text. To give a detailed explanation, we provide the pseudo codes in PyTorch style. Algorithm 1 and Algorithm 2 illustrate the first- and second-order methods, respectively. Moreover, we measure the disentanglement of five attributes by Re-scoring [31] in Fig. 8. The top row lists edited attributes, and the scores are the classification logits changes between original and edited images.

|  | Smile | Bangs | Gender | Glass | Age |
|---|---|---|---|---|---|
| Smile | 0.45 | -0.02 | 0.00 | -0.05 | -0.03 |
| Bangs | 0.00 | 0.52 | 0.00 | 0.00 | -0.01 |
| Gender | -0.03 | -0.03 | 0.54 | 0.02 | 0.03 |
| Glass | 0.00 | 0.00 | 0.01 | 0.52 | 0.01 |
| Age | -0.05 | -0.04 | 0.06 | 0.13 | 0.45 |

(a) Ours-1

|  | Smile | Bangs | Gender | Glass | Age |
|---|---|---|---|---|---|
| Smile | 0.45 | -0.01 | 0.00 | -0.03 | -0.02 |
| Bangs | 0.00 | 0.52 | 0.00 | 0.00 | 0.00 |
| Gender | -0.01 | -0.02 | 0.55 | 0.02 | 0.03 |
| Glass | 0.00 | 0.00 | 0.01 | 0.52 | 0.01 |
| Age | -0.02 | -0.03 | 0.04 | 0.12 | 0.45 |

(b) Ours-2

Figure 8. Re-scoring results of label-based editing on five attributes, Ours-1 and Ours-2 represent our first- and second-order methods, respectively.

## C. More Results

In this section, we provide more results of inversion, label-based editing and reference-based editing in Fig. 9, Fig. 10, Fig. 11.

---

**Algorithm 1** First-order Label-based Editing

```
1   # w: input latent code (18, 512)
2   # C: latent classifier
3   # y_t: target label
4
5   predicted = C(w)
6   loss = torch.nn.BCELoss(predicted, y_t)
7   loss.backward()
8   direct = w.grad
9   direct = direct / torch.norm(direct, dim=1)
10  w_edit = w - alpha * direct # alpha is a
        scaling factor.
```

---

**Algorithm 2** Second-order Label-based Editing

```
1   # w: input latent code (18, 512)
2   # C: latent classifier
3   # y_t: target label
4
5   r_d = torch.randn(18, 512)
6   r_0 = torch.zeros(18, 512)
7   w_d = w + kasi * r_d # kasi is a small number
        , we set it to 10e-4.
8   w_0 = w + r_0
9   predicted_d = C(w_d)
10  loss = torch.nn.BCELoss(predicted_d, y_t)
11  loss.backward()
12  direct_d = r_d.grad
13
14  C.zero_grad()
15  predicted_0 = C(w_0)
16  loss = torch.nn.BCELoss(predicted_0, y_t)
17  loss.backward
18  direct_0 = r_0.grad
19
20  direct = direct_d - direct_0
21  direct = direct / torch.norm(direct, dim=1)
22  w_edit = w - alpha * direct # alpha is a
        scaling factor.
```
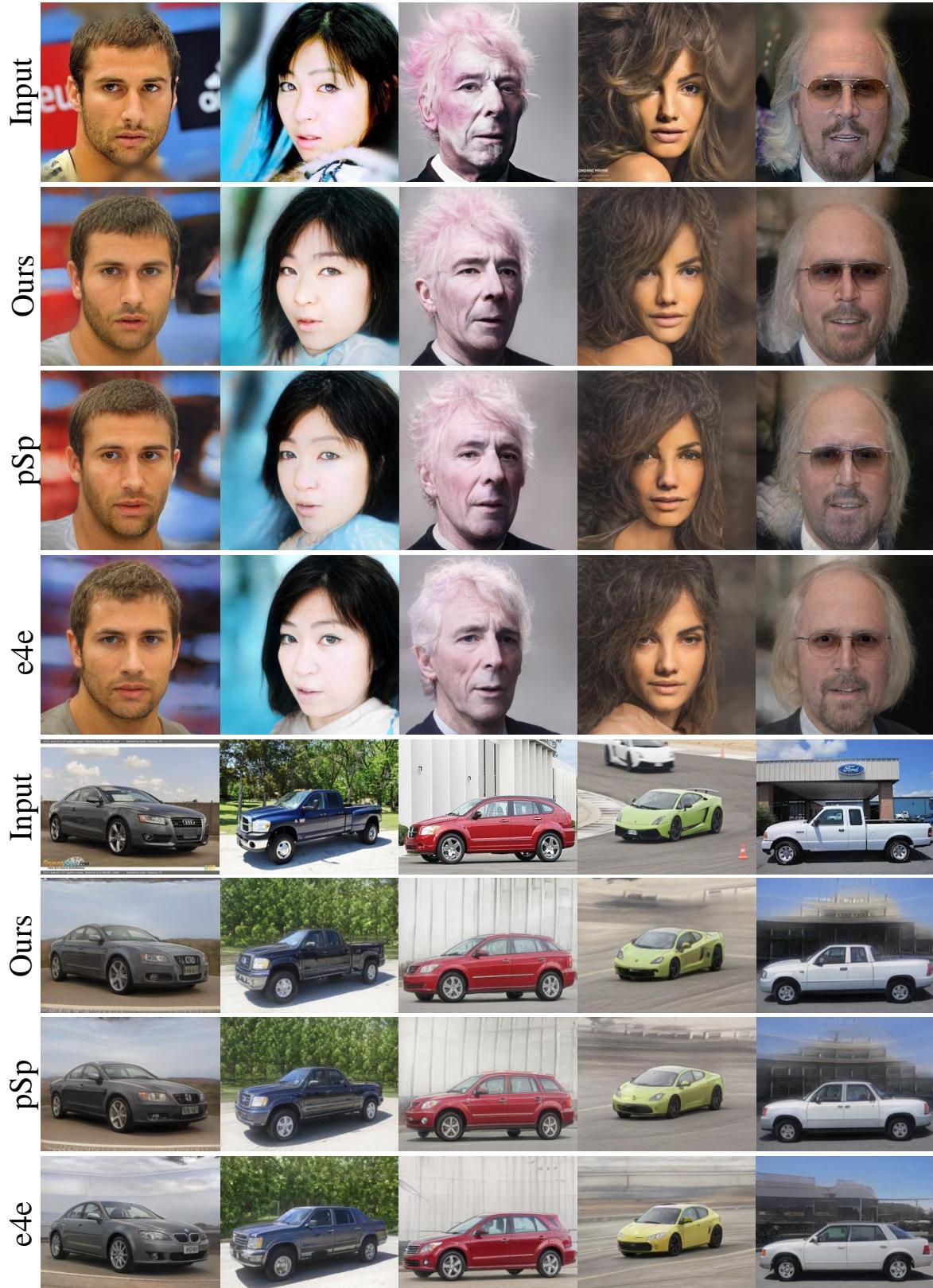
---

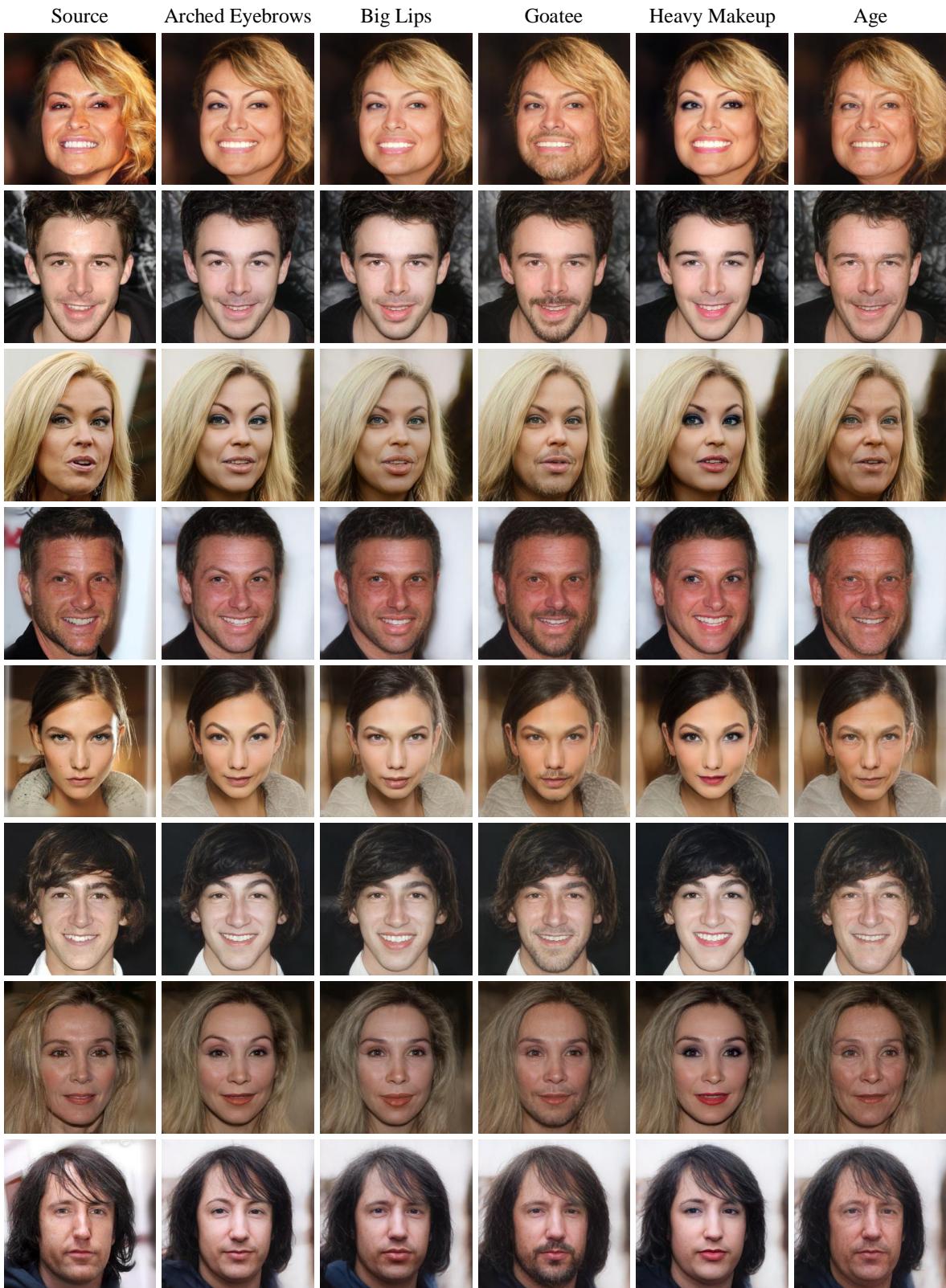Figure 9. More results of inversion compared with [30] and [33].

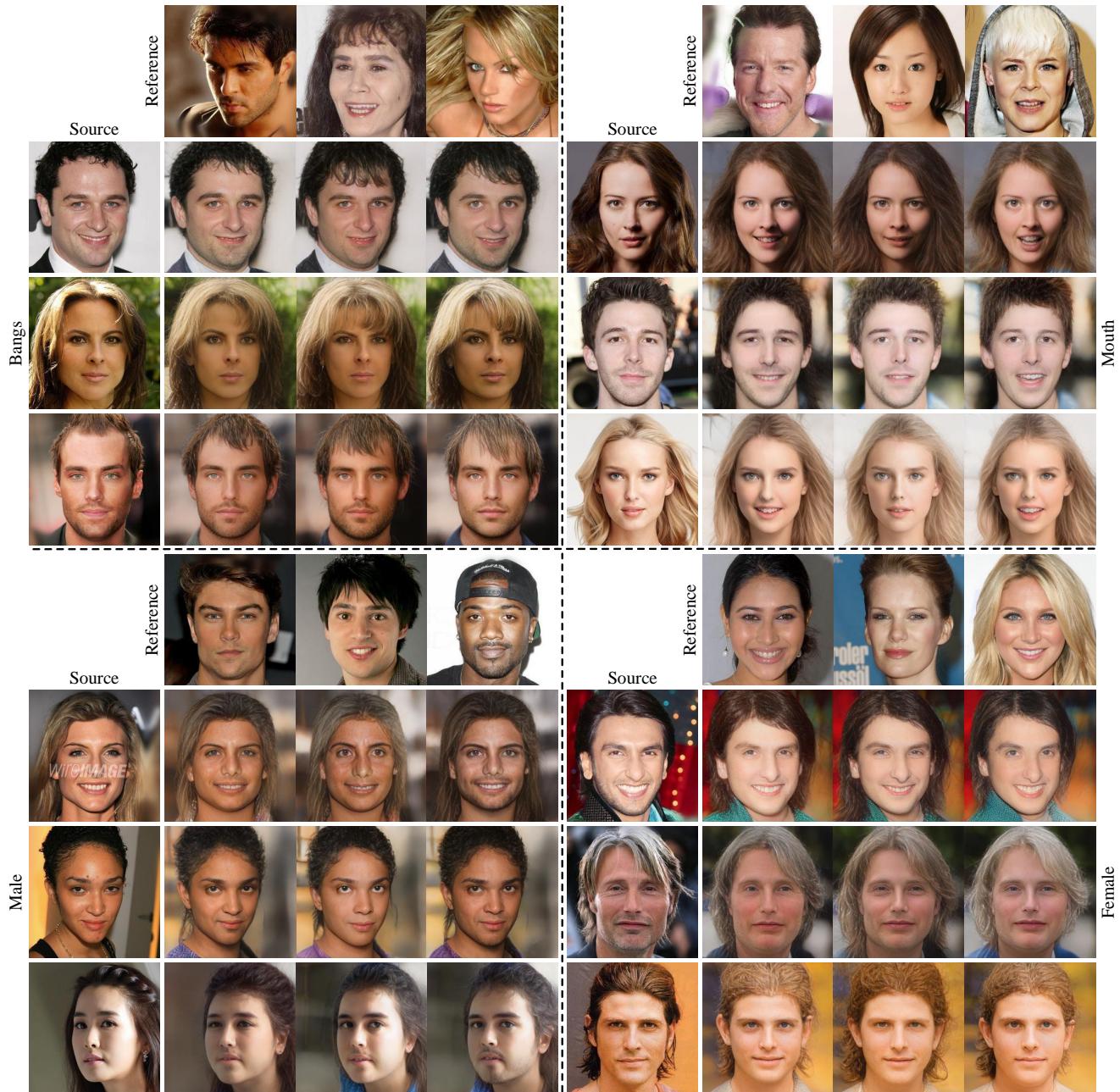Figure 10. More results of label-based editing on five attributes.

Figure 11. More results of reference-based editing on three attributes. The edited images take the style of *Bangs*, *Mouth* and *Gender* from different reference images.