

FEAT: Face Editing with Attention

Xianxu Hou¹, Linlin Shen^{1*}, Or Patashnik², Daniel Cohen-Or², Hui Huang¹
¹Shenzhen University ²Tel Aviv University



Figure 1. Manipulating local facial attributes based on text descriptions. The input and edited images are shown in the first and second row, respectively. The corresponding text prompts are shown at the bottom.

Abstract

Employing the latent space of pretrained generators has recently been shown to be an effective means for GAN-based face manipulation. The success of this approach heavily relies on the innate disentanglement of the latent space axes of the generator. However, face manipulation often intends to affect local regions only, while common generators do not tend to have the necessary spatial disentanglement. In this paper, we build on the StyleGAN generator, and present a method that explicitly encourages face manipulation to focus on the intended regions by incorporating learned attention maps. During the generation of the edited image, the attention map serves as a mask that guides a blending between the original features and the modified ones. The guidance for the latent space edits is achieved by employing CLIP, which has recently been shown to be effective for text-driven edits. We perform extensive experiments

*Corresponding author

and show that our method can perform disentangled and controllable face manipulations based on text descriptions by attending to the relevant regions only. Both qualitative and quantitative experimental results demonstrate the superiority of our method for facial region editing over alternative methods.

1. Introduction

Recent years have witnessed the significant progress of Generative Adversarial Networks (GANs) [13]. Specifically, StyleGAN [21–24], one of the most celebrated GAN frameworks, can produce high fidelity human face images with unmatched photo-realism. Furthermore, it has been shown that StyleGAN provides semantically rich latent space, which allows us to edit images in a semantically meaningful manner.

Numerous GAN-editing works use direction in the latent space to edit an image [14, 36, 37, 41]. However, using such

a direction for manipulation heavily relies on the assumption that the latent space is perfectly disentangled. Furthermore, these works typically require manual tuning of hyperparameters such as manipulation strength and disentanglement magnitude. Other works take an alternative approach and train a network that predicts a per-image offset in the latent space for an intended edit [2, 3, 17, 29]. These methods neither assume perfect disentanglement in the latent space, nor require manual tuning of the edit magnitude.

The most evident challenge in the aforementioned works is how to achieve a disentangled edit, that is, editing the intended attribute without affecting other attributes. In many cases, the disentanglement is achieved naturally, without explicit enforcement. However, some cases are more challenging and require special means. Previous works achieve disentanglement by constructing a spatially-disentangled representation of the latent space [19], or by operating in other latent spaces, for example, StyleSpace [11, 41, 46]. However, these methods heavily rely on the innate disentanglement of the latent space axes, which do not always have a spatial disentanglement.

In this paper, we present a method that explicitly encourages the edit to be focused on the intended region by incorporating the *learned* attention map. Specifically, given a latent code and a target edit, we train a network that predicts an offset in the latent space. During training, we learn an attention map by accounting for the features of all layers. The attention module that computes the masks is also guided by the edited attribute. The attention map is applied on the features of a target layer during the generation of the edited image. Blending the masked features with the original ones leads to the manipulation of the intended regions only. To guide the edit, we use CLIP [31], which has been shown to be effective for a wide range of image manipulations, and thus enables a simple and intuitive user interface (see Fig. 1 and Fig. 2).

To validate the effectiveness of the attention maps, we perform extensive experiments and show that we are able to obtain image changes in the intended regions by only using text descriptions. Moreover, we analyze the choice of the layers in which the attention maps are applied and show that it is necessary to apply them in the feature space. Finally, we compare our method with previous methods and show the advantage of using such attention maps.

2. Related Work

GAN Latent Space Manipulation. Generative Adversarial Networks have demonstrated great potential in learning various high-level semantics from observed data in the latent space, which provides an intuitive way to control the image generation process with different semantic specifications. An active line of research on image editing is to directly manipulate the latent code of a pretrained GAN

model. Early works [16, 32] observed interpretable vector arithmetic, which enables easy editing of different visual concepts in the learned latent space. Most notably, the discovery of such interpretable directions have received much research attention with the advancement of StyleGAN. Recent works have sought to identify a semantically meaningful path in a supervised manner, which requires a large number of annotated images [2] or attribute predictors on the predefined semantics [17, 36].

Two recent works [14, 37] have considered unsupervised identification of interpretable controls over image synthesis by using a closed-form factorization method or applying PCA in the latent space. While these approaches do not require supervision, they can only find limited semantic directions, which are selected manually and with extensive effort. Recently, domain adaptation methods [18, 38, 56] that build upon StyleGAN have also been proposed to achieve style transfer. Several works [11, 19, 25] have focused on local semantically-aware edits to a target output image by encoding the local semantics of images into the StyleGAN latent space. Additionally, in order to support real image manipulations, GAN inversion [1, 4, 5, 33, 34, 39, 43, 54, 57] has been adopted to inversely project a given image into the latent space of a pretrained GAN generator, and then the obtained codes can be edited in a semantically meaningful manner.

Text-based Image Synthesis and Manipulation. GANs have also been used to synthesize natural images based on text descriptions. For example, in some works, a conditional GAN has been used to generate bird and flower images by explicitly taking a language embedding as input. StackGAN [51] and StackGAN++ [52] use multiple generators and discriminators to further improve the quality of generated images. AttnGAN [45] uses an attentional generative network to select the condition at the word level for synthesizing fine-grained details in different subregions of the image. The semantic consistency and image fidelity can be further improved through carefully designed network architecture and training objectives [10, 26, 30, 49, 55]. More recently, Gal et al. [12] present a text-based zero-shot domain-adaptation technique.

Beyond text-to-image generation, manipulating face images using text descriptions has recently become popular with the advancement of StyleGAN. TediGAN [42] achieves text-based image generation and manipulation by mapping the image and text descriptions into a common StyleGAN latent space. StyleCLIP [29] uses Contrastive Language-Image Pretraining (CLIP) model [31] to develop a text-based interface for StyleGAN image manipulation. While these approaches can provide a certain level of text-based global attribute control, the local linguistic information of key words is not really attended. As a result, un-

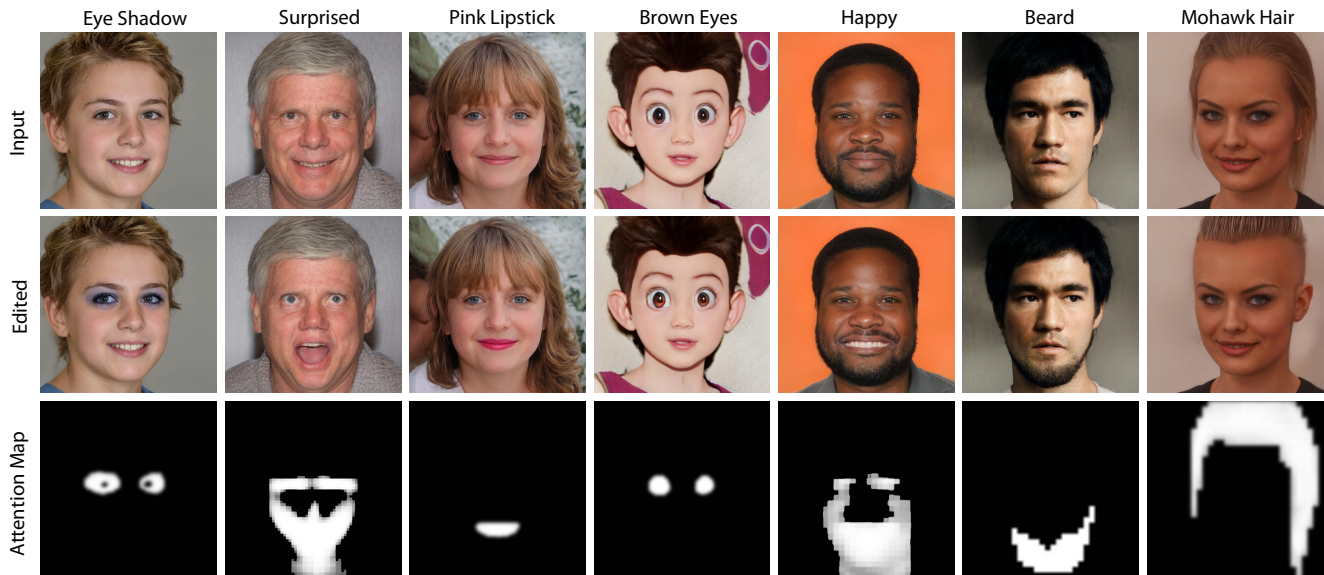


Figure 2. Edited results of our methods using different text prompts. The input and edited results are shown in the first and second row, respectively. The corresponding attention masks are shown in the third row.

wanted changes or artifacts may be produced. Our goal is to manipulate the semantics of a given image by automatically attending to the relevant regions.

Attention Learning. Attention learning has been successfully introduced in many applications in natural language processing and computer vision, such as image captioning [44], object localization [28], visual question answering [48] and machine translation [7]. Learning attention can also encourage more realistic image generation and manipulation. For example, SAGAN [50] considers the non-local relationships in the feature space by incorporating a self-attention module into the GAN framework. Moreover, several unsupervised image-to-image translation techniques [6,9,27,47] use attention network to predict spatial attention maps of images, demonstrating significant improvements in the quality of the translated images. A recently proposed method [8] enables users to manipulate a semantics of a GAN generated image at a given location, however it requires to manually specify the locations. In addition, a proposed facial structure editing method [40] aims to remove the double chin in portrait images relying on neck masks, which are extracted by a pretrained facial parsing model. By contrast, in our work we achieve controllable face manipulation by automatically attending to the regions of interest using only full-text descriptions.

3. Method

3.1. Preliminaries

StyleGAN2 is the current state of the art model for high-resolution image generation owing to its unique generator architecture. Instead of directly feeding the input latent code $z \in \mathcal{Z}$ to the network, it uses a mapping network f to transform z to the intermediate latent code $w = f(z)$, which is then transformed to style codes that control the layers of synthesis network g by modulating each layer’s convolutions. In layer $i + 1$ of g , the modulated convolution is applied on the features of layer i , denoted by $S_i(w_i)$. The output image is obtained by transforming StyleGAN2’s features to RGB features via toRGB layers. StyleGAN2’s intermediate latent space, \mathcal{W} , has been shown to be semantically meaningful and disentangled. These properties make StyleGAN2 very effective for image manipulation. It has become a common practice to edit images by traversing in StyleGAN2’s latent space.

3.2. Overview

Our method, FEAT, employs StyleGAN2’s latent space to edit an image (see Fig. 3). FEAT is built on previous methods that train a mapper h that learns an offset in the latent space to edit the image. To attain a disentangled edit, and to edit the intended regions only, FEAT learns an attention map to blend the features obtained by the source latent code with the features of the shifted latent code. To guide the edit, we employ CLIP [31], which allows for using text to learn the offset and generate the attention map.

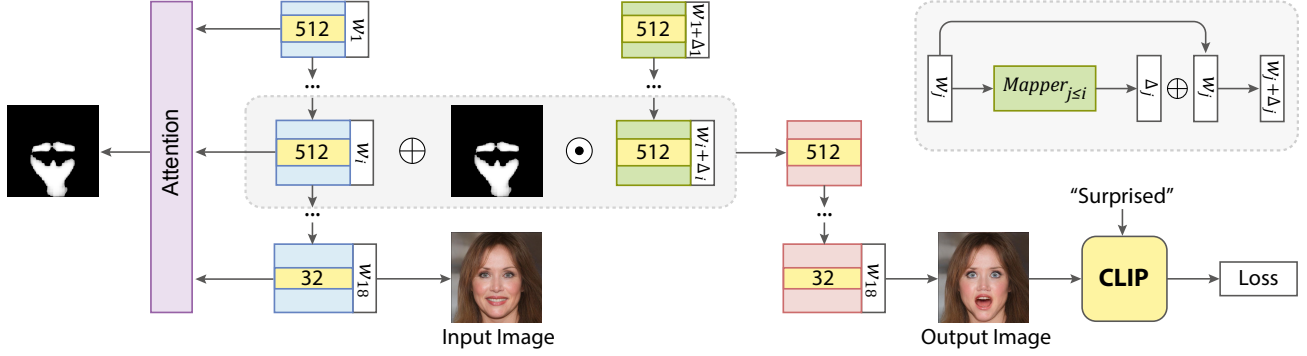


Figure 3. Overview : given an image generated by $w = (w_1, \dots, w_{18}) \in \mathcal{W}^+$ (the features are in light blue, where the number of channels are marked in the yellow band) and guided by a text prompt, we train mappers for style codes w_j for all $j \leq i$ that predict corresponding offsets Δ_j . We refer to the image generated by the modified $w_j + \Delta_j$ as the mapped image. In addition, we train an attention module (in light purple) that combines the features of all layers into a single attention map (on the left). Then, the i^{th} layer features of the original image and the mapped one are blended using the learned attention map.

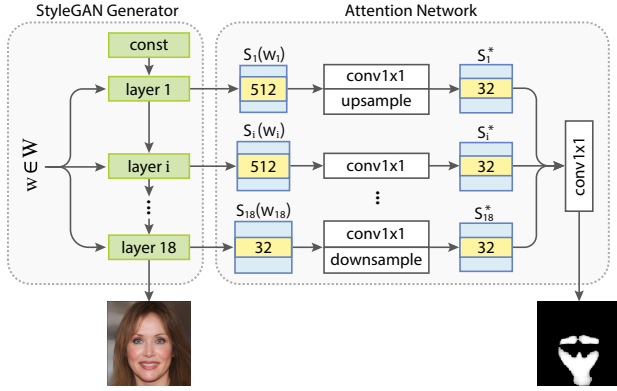


Figure 4. The architecture of our attention network f_{att} . The left part is the fixed StyleGAN2 generator for feature extraction. The right part is the attention network used to produce the attention masks.

3.3. The Mapper and Attention Module

Our mapper network h is implemented as an MLP that receives a latent code w and outputs an offset. The edited code is obtained by predicting the residual from the original latent code:

$$w_{edited} = w + \alpha \cdot h(w), \quad (1)$$

where α is a hyper-parameter used to balance between the original and edited data.

To edit the intended regions only, we employ an attention network f_{att} that learns to produce an attention map m , in which each pixel has a probability value between 0-1. The attention map is applied on the i^{th} layer features, $S_i(w_{edited})$, obtained during the generation of the image corresponding to w_{edited} . The attention map has a single channel, and the spatial dimensions of S_i . We apply m to

the features via an element-wise product in each channel to create the attended area: $m \odot S_i(w_{edited})$ (green features in Fig. 3). Similarly, the unattended area is created by $(1 - m) \odot S_i(w)$ (blue features). Thus, the final features (red features) are defined as follows:

$$\tilde{S}_i = m \odot S_i(w_{edited}) + (1 - m) \odot S_i(w). \quad (2)$$

Having obtained these blended features, the final edited image is attained by applying layers $> i$ on \tilde{S}_i , with w modulating the corresponding convolutions. During testing, we threshold the learned attention maps by only containing pixels exceeding a user-defined threshold τ , which we set to 0.8 in our experiments.

Inspired by the recent work [53], we exploit the features of all StyleGAN2 layers for the training of the attention network f_{att} . Fig. 4 illustrates the attention network architecture. We first use multiple 1×1 convolutional layers to reduce the number of channels of the feature maps $S_1(w_1), \dots, S_i(w_i), \dots, S_{18}(w_{18})$. For simplicity, all the feature maps are reduced to the same number of channels (32 by default in our implementation). The resulting feature maps are resized to the resolution of $S_i(w_i)$ and concatenated to a single deep feature map $S^* = [S_1^*, \dots, S_i^*, \dots, S_{18}^*]$. Here, S_i^* denotes the transformed feature map in the i^{th} layer, and S^* is then fed to another 1×1 convolution to produce a single channel feature map as the attention mask. A sigmoid layer is used at the end to ensure the output mask is bounded in the (0, 1) interval.

3.4. Training Objectives

Semantic Consistency Loss. First, we adopt a text-image matching network to define a semantic consistency loss, which guides the image manipulation to match the provided

text descriptions. In particular, given the text description t and the pretrained CLIP, our semantic consistency loss is defined as:

$$\mathcal{L}_{clip} = D_{clip}(I^*, t) \quad (3)$$

where D_{clip} denotes the cosine distance between the image and text embeddings extracted with CLIP, and I^* is the blended output image. Note that the features of the blended image outside the mask remain intact. Thus, the CLIP model is encouraged to focus the optimization on the attended area only, and the attention module encourages semantic consistency.

Attention Map Regularization. To encourage the attention network to focus on a more compact region related to the text descriptions rather than the whole image, we further introduce an attention map regularization as:

$$\mathcal{L}_{att} = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W m_{ij}, \quad (4)$$

where H and W are the width and height of the learned attention map, respectively, and m_{ij} is the pixel value at the location (i, j) .

Total Variation Loss. To encourage spatial smoothness in the produced attention map m , we adopt a total variation penalty as:

$$\mathcal{L}_{tv} = \sum_{i,j} \|m_{i+1,j} - m_{i,j}\|_2 + \|m_{i,j+1} - m_{i,j}\|_2. \quad (5)$$

Latent Loss. To explicitly preserve visual attributes of the input images, we minimize the L_2 distance between the original latent code w and the transformed code w_{edited} as:

$$\mathcal{L}_{l_2} = \|w - w_{edited}\|_2. \quad (6)$$

Overall Training Objective. Our full objective is the weighted sum of these four losses:

$$\mathcal{L} = \mathcal{L}_{clip} + \lambda_{att}\mathcal{L}_{att} + \lambda_{tv}\mathcal{L}_{tv} + \lambda_{l_2}\mathcal{L}_{l_2}, \quad (7)$$

where λ_{att} , λ_{tv} and λ_{l_2} are the hyper-parameters to control the relative importance of each component.

4. Results and Evaluation

4.1. Experimental Settings

Datasets. We use the StyleGAN2 model pretrained on FFHQ dataset [23] as our generator. The FFHQ dataset is a 1024×1024 resolution high-quality image dataset of human faces. We adopt e4e [39] to project the test set of CelebA-HQ [20] for real image manipulation. We also use

Table 1. Quantitative comparison with TediGAN and StyleCLIP by using different metrics. \downarrow indicates that lower is better. \uparrow indicates that higher is better.

Attribute	FID \downarrow			CS \uparrow			ED \downarrow		
	TG	SC	Ours	TG	SC	Ours	TG	SC	Ours
Curly hair	32.37	15.61	9.87	0.52	0.75	0.78	0.93	0.67	0.62
Mohawk hair	22.78	14.81	10.02	0.49	0.75	0.81	0.96	0.68	0.59
Purple hair	27.17	11.13	13.43	0.65	0.82	0.78	0.81	0.40	0.62
Surprised	11.73	9.31	4.73	0.68	0.89	0.96	0.78	0.43	0.27
Angry	14.46	6.15	2.51	0.59	0.87	0.88	0.88	0.49	0.47
Average	21.70	11.40	8.11	0.59	0.82	0.84	0.87	0.53	0.51

Notations: TG - TediGAN; SC - StyleCLIP; FID - Fréchet Inception Distance; CS - Cosine Similarity; ED - Euclidean Distance.

StyleGAN2 model fine-tuned on cartoon faces for evaluation (see, for example, the rightmost column in Fig. 1). In addition, we perform the editing on real faces collected from the Internet.

Baselines. TediGAN [42] and StyleCLIP [29] are closely related works, as both can achieve semantic face manipulation based on text descriptions. In addition, we use InterfaceGAN [36] and StyleFlow [2] as our baselines, which control face manipulations along predefined semantic directions.

Evaluation Metrics. We evaluate both the visual quality and identity preservation using quantitative metrics. Fréchet Inception Distance (FID) [15] is used to measure the discrepancy between the edited and original faces. We adopt the face recognition model FaceNet [35] to extract the embeddings of tested images and use Cosine Similarity (CS) and Euclidean Distance (ED) to quantify the identity preservation.

Training Details. We adopt a three-layer MLP with 512 hidden units and 512 output units to build our mapper network. Image sampling is performed by randomly drawing from a normal distribution from the \mathcal{Z} space, which is then mapped into the intermediate latent space \mathcal{W} . The output images are resized to 224×224 before feeding them to CLIP (using the pretrained ViT-B/32 weights). We use Adam optimizer to train the mapper and attention network at a learning rate of 0.0001. We use a batch size of 2 on one Tesla V100 32GB GPU. The maximum number of iterations is set to 20,000. Note that the pretrained StyleGAN2 is fixed. We set $\lambda_{att} = 0.005$, $\lambda_{tv} = 0.00001$, $\lambda_{l_2} = 0.8$, and $\alpha = 0.1$. We edit the first eight layers for facial structure manipulation such as hairstyle and expression, and edit all the 18 layers when performing local color manipulation, such as that for the hair and eyes.

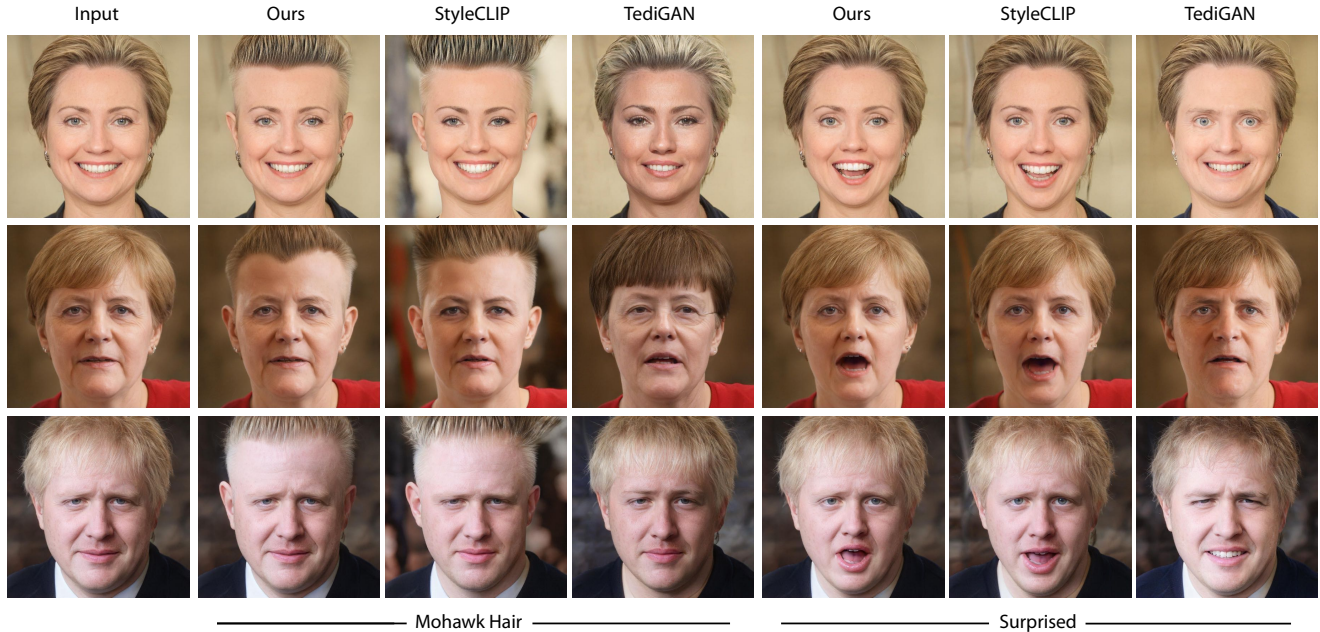


Figure 5. Visual comparison of our method with TediGAN and StyleCLIP. The driving descriptions are indicated at the bottom.

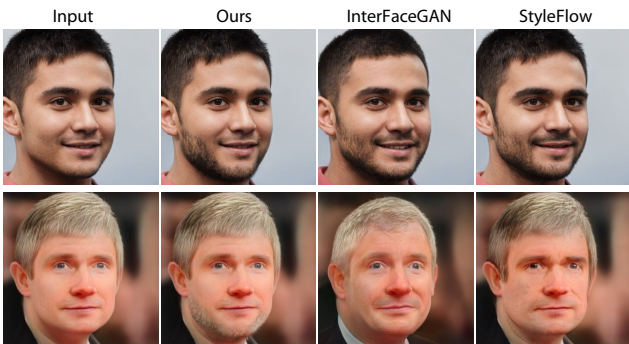


Figure 6. Visual comparison of our method on "beard" manipulation with InterFaceGAN and StyleFlow.

4.2. Learning to Attend

In this section, we inspect the attention module learned by our framework. Fig. 2 shows the produced attention maps when editing different facial attributes. We can observe that our approach succeeds in attending to the most relevant facial regions and creating semantic manipulation based on the text descriptions. We can see that in the example of eye shadow editing (the 1st column), the attention map has large values in the vicinity of the eyes, as expected. It can also be observed that the learned attention maps can correctly localize the eyes, lip, beard and hair, thus avoiding unnecessary changes to irrelevant areas when performing

editing.

4.3. Qualitative and Quantitative Comparisons

First, we compare our approach with two recently proposed text-based manipulation methods: TediGAN and StyleCLIP. TediGAN encodes both the image and the text into the StyleGAN latent space and can support text-based image manipulation with CLIP model. StyleCLIP investigates three techniques that combine CLIP with StyleGAN. We use the latent mapper approach, which is closer to our architecture, for comparison.

In addition, we provide comparative results with other prominent StyleGAN manipulation methods: InterFaceGAN and StyleFlow. InterFaceGAN performs linear manipulation in the GAN latent space, while StyleFlow extracts non-linear editing paths in the latent space using conditional continuous normalizing flows.

Qualitative Results. Fig. 5 presents the visual comparison with TediGAN and StyleCLIP by using different text descriptions. As can be observed, our method achieves more precise control over facial attributes, including the Mohawk hairstyle and the surprised expression. It can be seen that TediGAN fails to change the hairstyle and could alter the skin color. In addition, the perceived gender (the second example in the rightmost column) changes when performing expression manipulation. StyleCLIP can produce better results than TediGAN and successfully achieve

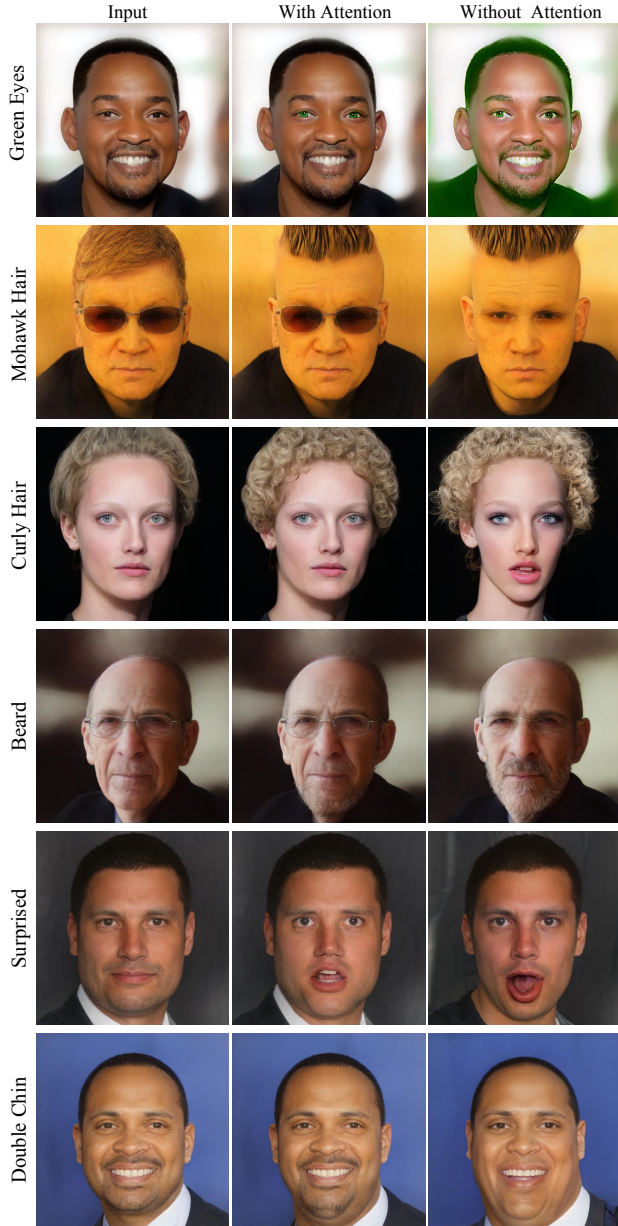


Figure 7. The effectiveness of the attention mechanism. We compare the edited results of our full model (the second column) and the one without the attention module (the third column). The corresponding text prompts are shown on the side.

these editing. However, other attributes can be heavily altered. We can observe that there exist noticeable changes of the background when editing the hairstyle and expression. By contrast, with the guidance of the attention map, our model effectively avoids unwanted changes by attending only to the target regions.

In Fig. 6, we provide a visual comparison with InterFaceGAN and StyleFlow. We consider adding a beard on the

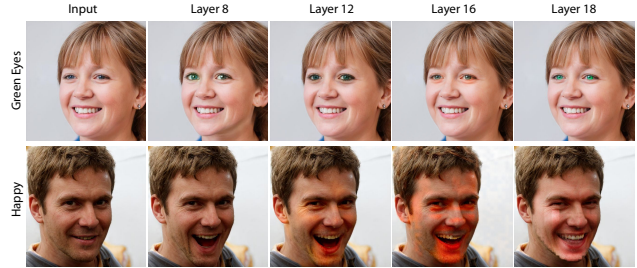


Figure 8. Visual comparison by performing blending in different feature layers. The corresponding text prompts are shown on the side. We can clearly tell that for color manipulation Layer 18 is the best, while for structural manipulation Layer 8 is the best.

faces, which the compared methods succeed in performing. Note that these two approaches require labeled data for supervision, and unlike our method, they cannot apply zero-shot manipulations. We notice that InterFaceGAN and StyleFlow may cause unwanted changes to the eyebrows when asked to add a beard. Compared to these approaches, our method is able to perform more controllable edits by only focusing on the most relevant parts in the images.

Quantitative Results. The superiority of our method can also be validated by quantitative evaluation. We evaluate FID, CS and ED based on 10,000 samples randomly generated with the StyleGAN2 generator. The results of five attributes editing compared with TediGAN and StyleCLIP are tabulated in Table 1. It can be seen that our method outperforms other approaches in terms of both visual quality (lower FID values) and identity preservation (higher CS and lower ED values). Our results are more favorable because our model can correctly localize the area of interest, thus avoiding unnecessary changes to irrelevant parts.

4.4. Attention Mechanism Analysis

Attention Module. Our attention module is used to produce a probability map that decides the regions to be modified. To demonstrate the effectiveness of the proposed attention mechanism, we perform an ablation study by muting the attention component. We show several examples in Fig. 7. As can be seen, without the guidance of the attention map, different attributes are entangled. In the example of facial structure manipulation, the eyeglasses and face shape (the 2nd and 3rd row) are entangled with the hairstyle. We can also see obvious changes of the background and clothes when editing the surprised expression (the 5th row). Moreover, the entanglement is particularly noticeable when performing color manipulation, and we can see that the background color and skin tone change significantly when we try to achieve green eyes (the 1st row). By contrast, with the

guidance of the attention map, our algorithm can effectively avoid unwanted changes by attending only to the target regions.

Blending in Different Layers. In this part, we explore the effect of feature blending in different layers. Two examples are shown in Fig. 8. As can be observed, our approach performs well on the color prompt (the 1st row) when performing the blending in the highest layer (Layer 18), and blending in the lower layer (Layer 8) results in changes of the eye shape. By contrast, for facial structural editing like the happy expression (the 2nd row), it is better to adopt lower layers for feature blending. This is because the fact that StyleGAN controls high-level aspects such as the pose, hair style and face shape in the coarse and middle layers ($4^2 - 32^2$) and deals with colors and smaller details in the finer layers ($64^2 - 1024^2$).

Two-Step Image Manipulation. In addition to simultaneously learning the attention masks and performing various face editing, we further explore a two-step training approach, where we use the first prompt to define the mask and the second one for manipulation with the mask produced in the first step. In this way, we are able to control facial regions with more interesting descriptions. Fig 9 shows four examples. In particular, we first use the “Black Hair” prompt to extract the hair area and fix the obtained mask to train models using the text prompt “Fire”, “Mosaic”, “Lavender” and “Willow Leaf”, respectively. In this experiment, we perform the feature blending at the 12th layer. As can be seen, by explicitly limiting the area that can be edited, our model can achieve more meaningful manipulation in the specific facial region. These results demonstrate that the attention mechanism can assist with an intuitive control over facial areas using text descriptions.

5. Conclusion and Future Work

We presented a novel method for semantic editing of facial regions guided by text descriptions. Our approach features a novel attention framework that leverages the knowledge learned from text-image joint embedding to guide the image generation process of StyleGAN. The attention module learns attention maps with text supervision only and is particularly useful to focus the manipulation on the intended regions of interest. Unlike previous methods that strongly rely on the disentanglement of different attributes in the latent space, our approach alleviates this dependency by explicitly restricting the altered spatial regions. Experiments demonstrate the superior performance of our method over previous works.

In the future, it will be interesting to explore the potential of these attention masks as means to segment facial regions.

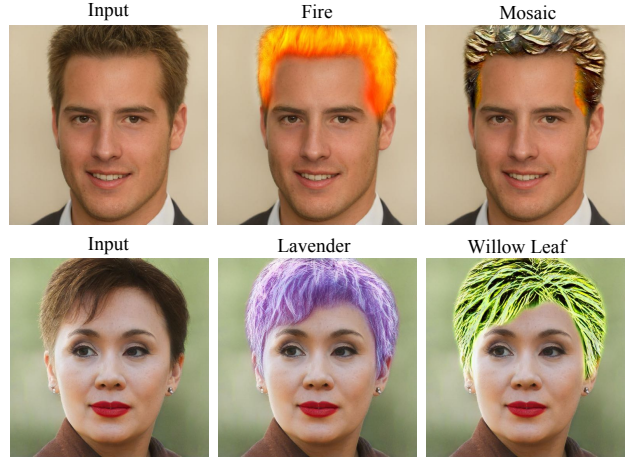


Figure 9. Two-step image manipulation. In the first step, we train the attention module with the text prompt “Black Hair”. In the second step, the attention map is fixed, and the mapper is trained with the text prompts “Fire”, “Mosaic”, “Lavender” and “Willow Leaf”, respectively.

Using guiding text or possibly a few-shot, our mapper can detect the corresponding facial regions. Such segmentation may successfully be generalized to various faces with different poses. Furthermore, we believe that such attention mechanisms can be extended beyond the facial domain and bypass attributes that are spatially entangled.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proc. Int. Conf. on Computer Vision*, pages 4432–4441, 2019. 2
- [2] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *ACM Trans. on Graphics*, 40(3):21:1–21:21, 2021. 2, 5
- [3] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Only a matter of style: Age transformation using a style-based regression model. *ACM Trans. on Graphics (Proc. SIGGRAPH)*, 40(4):45:1–45:12, 2021. 2
- [4] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder via iterative refinement. In *Proc. Int. Conf. on Computer Vision*, pages 6711–6720, 2021. 2
- [5] Yuval Alaluf, Omer Tov, Ron Mokady, Rinon Gal, and Amit H. Bermano. Hyperstyle: Stylegan inversion with hypernetworks for real image editing. *arXiv preprint arXiv:2111.15666*, 2021. 2
- [6] Youssef Alami Mejjati, Christian Richardt, James Tompkin, Darren Cosker, and Kwang In Kim. Unsupervised attention-guided image-to-image translation. In *Proc. Conf. on Neural Information Processing Systems*, pages 3697–3707, 2018. 3

- [7] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014. **3**
- [8] David Bau, Alex Andonian, Audrey Cui, YeonHwan Park, Ali Jahanian, Aude Oliva, and Antonio Torralba. Paint by word. *arXiv preprint arXiv:2103.10951*, 2021. **3**
- [9] Xinyuan Chen, Chang Xu, Xiaokang Yang, and Dacheng Tao. Attention-gan for object transfiguration in wild images. In *Proc. Euro. Conf. on Computer Vision*, pages 164–180, 2018. **3**
- [10] Jun Cheng, Fuxiang Wu, Yanling Tian, Lei Wang, and Dapeng Tao. Rifegan: Rich feature generation for text-to-image synthesis from prior knowledge. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 10911–10920, 2020. **2**
- [11] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 5771–5780, 2020. **2**
- [12] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *arXiv preprint arXiv:2108.00946*, 2021. **2**
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proc. Conf. on Neural Information Processing Systems*, pages 2672–2680, 2014. **1**
- [14] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. In *Proc. Conf. on Neural Information Processing Systems*, 2020. **1, 2**
- [15] 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. In *Proc. Conf. on Neural Information Processing Systems*, pages 6626–6637, 2017. **5**
- [16] Xianxu Hou, Linlin Shen, Ke Sun, and Guoping Qiu. Deep feature consistent variational autoencoder. In *Proc. IEEE Winter Conf. on Applications of Computer Vision*, pages 1133–1141, 2017. **2**
- [17] Xianxu Hou, Xiaokang Zhang, Hanbang Liang, Linlin Shen, Zhihui Lai, and Jun Wan. Guidedstyle: Attribute knowledge guided style manipulation for semantic face editing. *Neural Networks*, 145:209–220, 2022. **2**
- [18] Wonjong Jang, Gwangjin Ju, Yuchool Jung, Jiaolong Yang, Xin Tong, and Seungyong Lee. Stylecarigan: caricature generation via stylegan feature map modulation. *ACM Trans. on Graphics (Proc. SIGGRAPH)*, 40(4):116:1–116:16, 2021. **2**
- [19] Omer Kafri, Or Patashnik, Yuval Alaluf, and Daniel Cohen-Or. Stylefusion: A generative model for disentangling spatial segments. *arXiv preprint arXiv:2107.07437*, 2021. **2**
- [20] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *Proc. Int. Conf. on Learning Representations*, 2018. **5**
- [21] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Proc. Conf. on Neural Information Processing Systems*, 2020. **1**
- [22] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Proc. Conf. on Neural Information Processing Systems*, 2021. **1**
- [23] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 4401–4410, 2019. **1, 5**
- [24] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 8110–8119, 2020. **1**
- [25] Hyunsu Kim, Yunje Choi, Junho Kim, Sungjoo Yoo, and Youngjung Uh. Exploiting spatial dimensions of latent in gan for real-time image editing. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 852–861, 2021. **2**
- [26] Wenbo Li, Pengchuan Zhang, Lei Zhang, Qiuyuan Huang, Xiaodong He, Siwei Lyu, and Jianfeng Gao. Object-driven text-to-image synthesis via adversarial training. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 12174–12182, 2019. **2**
- [27] Ron Mokady, Sagie Benaim, Lior Wolf, and Amit Bermano. Masked based unsupervised content transfer. In *Proc. Int. Conf. on Learning Representations*, 2019. **3**
- [28] Maxime Oquab, Léon Bottou, Ivan Laptev, and Josef Sivic. Is object localization for free?-weakly-supervised learning with convolutional neural networks. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 685–694, 2015. **3**
- [29] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proc. Int. Conf. on Computer Vision*, pages 2085–2094, 2021. **2, 5**
- [30] Tingting Qiao, Jing Zhang, Duanqing Xu, and Dacheng Tao. Mirrorgan: Learning text-to-image generation by redescription. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 1505–1514, 2019. **2**
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proc. Int. Conf. on Machine Learning*, pages 8748–8763, 2021. **2, 3**
- [32] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. **2**
- [33] 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. *arXiv preprint arXiv:2008.00951*, 2020. **2**

- [34] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *arXiv preprint arXiv:2106.05744*, 2021. [2](#)
- [35] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 815–823, 2015. [5](#)
- [36] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 9243–9252, 2020. [1](#), [2](#), [5](#)
- [37] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 1532–1540, 2021. [1](#), [2](#)
- [38] Guoxian Song, Linjie Luo, Jing Liu, Wan-Chun Ma, Chunpong Lai, Chuanxia Zheng, and Tat-Jen Cham. Agilegan: stylizing portraits by inversion-consistent transfer learning. *ACM Trans. on Graphics (Proc. SIGGRAPH)*, 40(4):117:1–117:13, 2021. [2](#)
- [39] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *ACM Trans. on Graphics (Proc. SIGGRAPH)*, 40(4):133:1–133:14, 2021. [2](#), [5](#)
- [40] Yiqian Wu, Yong-Liang Yang, Qinjie Xiao, and Xiaogang Jin. Coarse-to-fine: facial structure editing of portrait images via latent space classifications. *ACM Trans. on Graphics (Proc. SIGGRAPH)*, 40(4):46:1–46:13, 2021. [3](#)
- [41] Zongze Wu, Dani Lischinski, and Eli Shechtman. Stylespace analysis: Disentangled controls for stylegan image generation. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 12863–12872, 2021. [1](#), [2](#)
- [42] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 2256–2265, 2021. [2](#), [5](#)
- [43] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. Gan inversion: A survey. *arXiv preprint arXiv:2101.05278*, 2021. [2](#)
- [44] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *Proc. Int. Conf. on Machine Learning*, pages 2048–2057, 2015. [3](#)
- [45] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 1316–1324, 2018. [2](#)
- [46] Yinghao Xu, Yujun Shen, Jiapeng Zhu, Ceyuan Yang, and Bolei Zhou. Generative hierarchical features from synthesizing images. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 4432–4442, 2021. [2](#)
- [47] Chao Yang, Taehwan Kim, Ruizhe Wang, Hao Peng, and C-C Jay Kuo. Show, attend, and translate: Unsupervised image translation with self-regularization and attention. *IEEE Trans. on Image Processing*, 28(10):4845–4856, 2019. [3](#)
- [48] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. Stacked attention networks for image question answering. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 21–29, 2016. [3](#)
- [49] Guojun Yin, Bin Liu, Lu Sheng, Nenghai Yu, Xiaogang Wang, and Jing Shao. Semantics disentangling for text-to-image generation. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 2327–2336, 2019. [2](#)
- [50] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *Proc. Int. Conf. on Machine Learning*, pages 7354–7363, 2019. [3](#)
- [51] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proc. Int. Conf. on Computer Vision*, pages 5907–5915, 2017. [2](#)
- [52] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE Trans. Pattern Analysis & Machine Intelligence*, 41(8):1947–1962, 2018. [2](#)
- [53] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 10145–10155, 2021. [4](#)
- [54] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image editing. In *Proc. Euro. Conf. on Computer Vision*, pages 592–608, 2020. [2](#)
- [55] Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. Dm-gan: Dynamic memory generative adversarial networks for text-to-image synthesis. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 5802–5810, 2019. [2](#)
- [56] Peihao Zhu, Rameen Abdal, John Femiani, and Peter Wonka. Mind the gap: Domain gap control for single shot domain adaptation for generative adversarial networks. *arXiv preprint arXiv:2110.08398*, 2021. [2](#)
- [57] Peihao Zhu, Rameen Abdal, Yipeng Qin, John Femiani, and Peter Wonka. Improved stylegan embedding: Where are the good latents? *arXiv preprint arXiv:2012.09036*, 2021. [2](#)

FEAT: Face Editing with Attention (Supplementary Material)

Xianxu Hou¹, Linlin Shen¹, Or Patashnik², Daniel Cohen-Or², Hui Huang¹
¹Shenzhen University ²Tel Aviv University

1. Qualitative Comparisons

Here, we provide additional comparisons of our method with the current state-of-the-art methods for text-based manipulation: TediGAN [5] and StyleCLIP [4]. These methods also combine CLIP with StyleGAN [2, 3] to achieve text-based editing. For StyleCLIP, we use the latent mapper approach for comparison, which is closer to our architecture. For structural manipulation, we also compare our method with StyleFusion [1], which provides fine-grained control over each region of the generated image by fusing different latent codes. In particular, we first use StyleCLIP to manipulate the latent code to obtain the edited representation, which is then fed to StyleFusion hierarchy to control the areas relevant to the desired edit. Then, we use the original latent code to control the remaining facial areas.

Fig. 1 presents the visual comparison using different color prompts. As can be seen, our method can achieve more precise control over various facial attributes, including different colors of eyes, hair, lipstick and eye shadow. We observe that TediGAN fails to produce colorful eyes and hair, and the perceived identity is changed when editing lipstick. StyleCLIP can produce better results than TediGAN and successfully change the colors of different facial attributes. However, the color of the background and clothes may also be altered when producing the green eyes.

Figs. 2, 3, 4 and 5 present the comparisons on several structural manipulations, including different hairstyles and facial expressions. It can be seen that TediGAN struggles in producing the Mohawk hairstyle and the perceived gender has been changed (Fig. 2). Although StyleCLIP can successfully achieve these edits, background can be entangled with hairstyle (Fig. 2) and expression (Fig. 4). StyleFusion can better preserve the background than StyleCLIP, however we can observe the tone change of the entire image (Fig. 4) and slight entanglement between the hair and eyes. By contrast, our approach is able to effectively avoid unwanted changes by providing precise local control over different attributes. Moreover, in StyleFusion the regions are pre-defined independently to the text-prompt, while in our method the regions are learned for the specific text-prompt.

References

- [1] Omer Kafri, Or Patashnik, Yuval Alaluf, and Daniel Cohen-Or. Stylefusion: A generative model for disentangling spatial segments. *arXiv preprint arXiv:2107.07437*, 2021. 1
- [2] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 4401–4410, 2019. 1
- [3] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 8110–8119, 2020. 1
- [4] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proc. Int. Conf. on Computer Vision*, pages 2085–2094, 2021. 1
- [5] Weihao Xia, Yujia Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, pages 2256–2265, 2021. 1



Figure 1. Visual comparison of our method with TediGAN and StyleCLIP. The driving descriptions are indicated above each column.

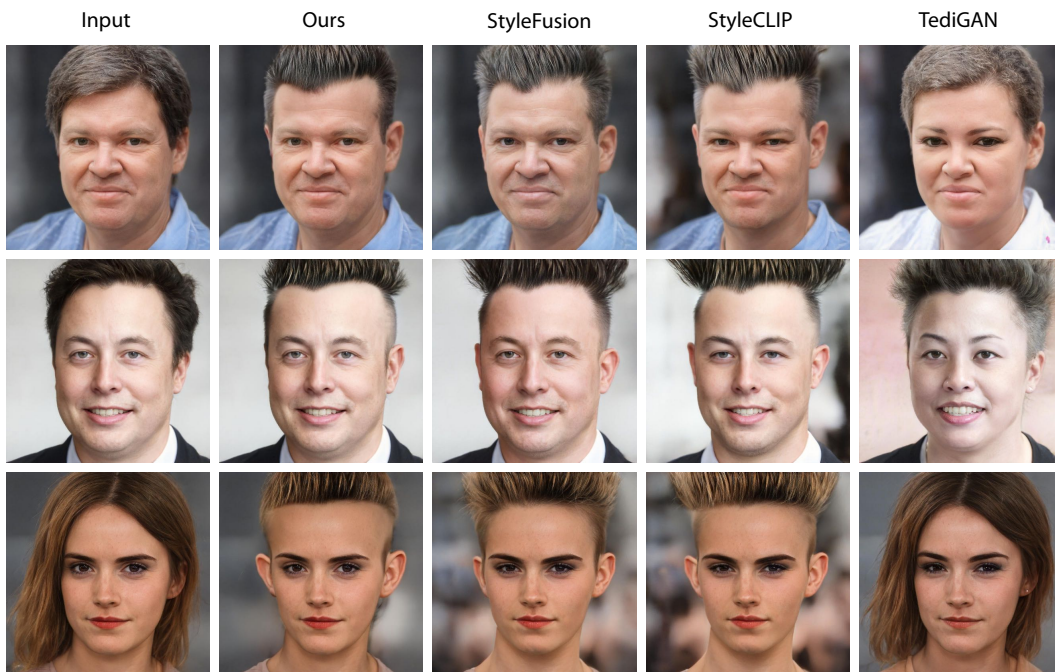


Figure 2. Visual comparison of our method with TediGAN, StyleCLIP and StyleFusion on Mohawk hairstyle. As can be seen, our method yields a disentangled edit. For example, in the last row the background is unnecessarily modified even when combining StyleCLIP with StyleFusion.

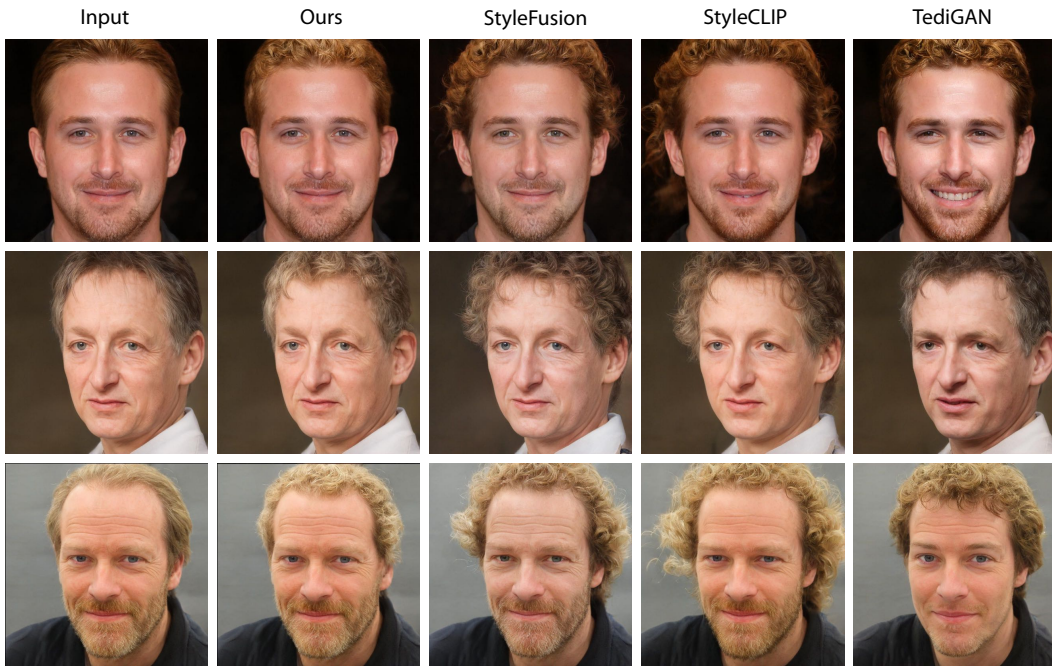


Figure 3. Visual comparison of our method with TediGAN, StyleCLIP and StyleFusion on curly hairstyle. As can be seen, while in other methods the hair style (curly) is entangled with the hair length, our method successfully alters only the hair style.

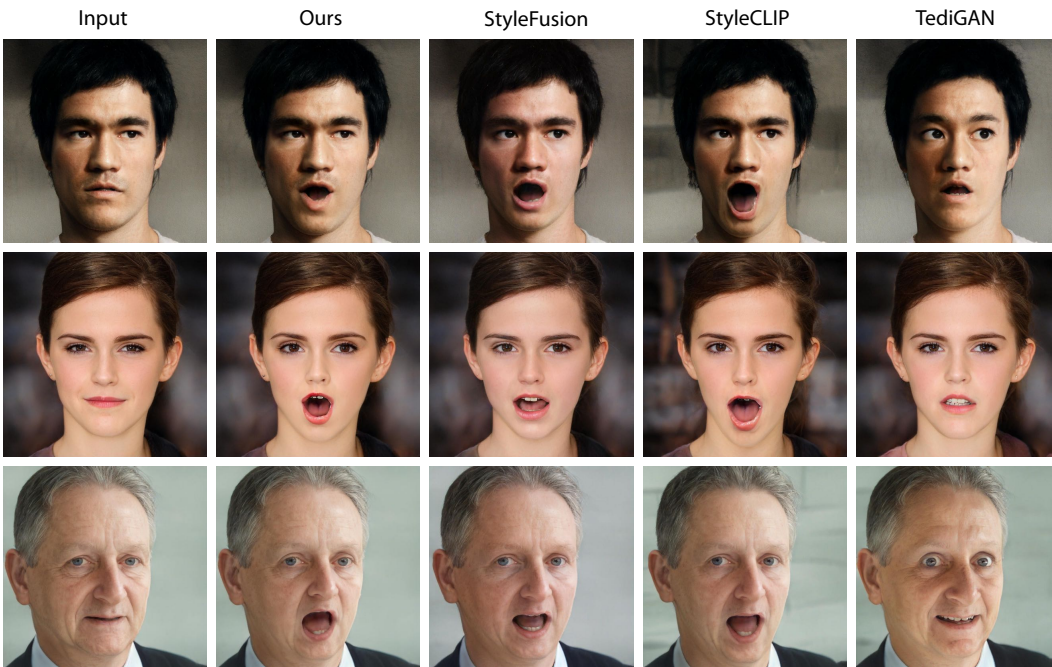


Figure 4. Visual comparison of our method with TediGAN, StyleCLIP and StyleFusion on surprised expression. As can be seen, our method achieves more faithful results compared to other works. For example, our method is more disentangled than StyleCLIP, which slightly changes the identity, and more accurate than StyleFusion, which does not alter the eyes.



Figure 5. Visual comparison of our method with TediGAN, StyleCLIP and StyleFusion on angry expression. For StyleFusion, we only edit the relevant attributes: face, eyes and mouth.