


Pastiche Master: Exemplar-Based High-Resolution Portrait Style Transfer

Shuai Yang Liming Jiang Ziwei Liu Chen Change Loy 

S-Lab, Nanyang Technological University

{shuai.yang, liming002, ziwei.liu, ccloy}@ntu.edu.sg

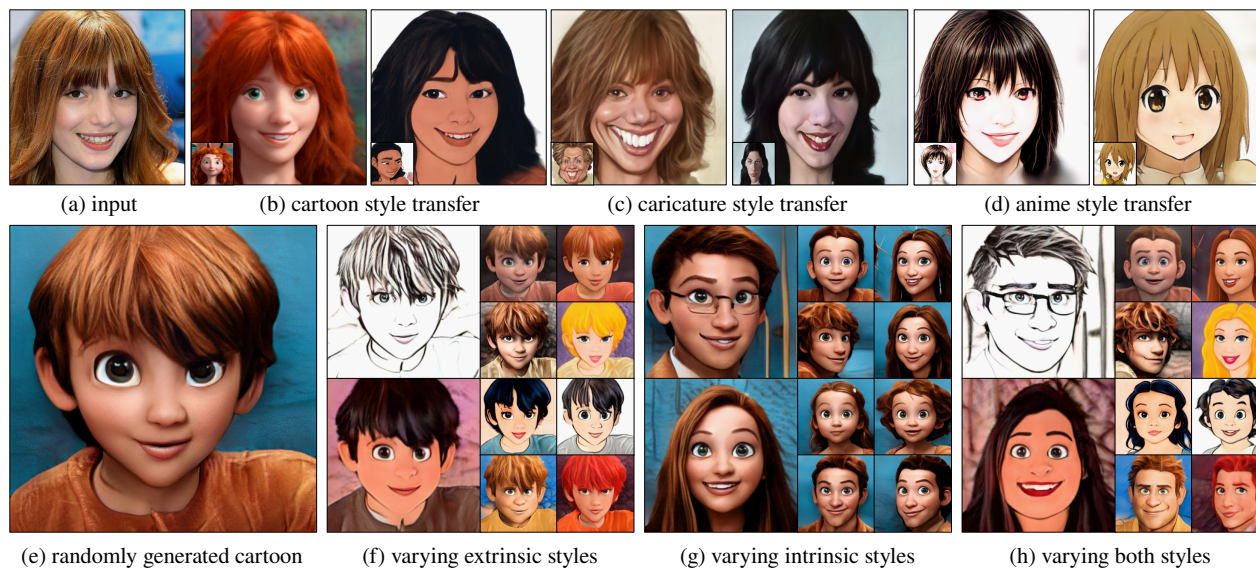


Figure 1. We propose a novel DualStyleGAN for exemplar-based high-resolution (1024×1024) portrait style transfer. The artistic portraits of (b)-(d) generated from a real face (a) successfully imitate the color and structural styles of the examples seen in their respective lower-left corners. DualStyleGAN features dual style paths: an intrinsic style path and an extrinsic style path for flexible control over the content and style, respectively. (e) A cartoon face generated from arbitrary intrinsic and extrinsic style codes. Samples generated by (f) varying extrinsic styles with fixed intrinsic styles, (g) varying intrinsic styles with fixed extrinsic styles, and (h) varying both styles.

Abstract

Recent studies on StyleGAN show high performance on artistic portrait generation by transfer learning with limited data. In this paper, we explore more challenging exemplar-based high-resolution portrait style transfer by introducing a novel **DualStyleGAN** with flexible control of dual styles of the original face domain and the extended artistic portrait domain. Different from StyleGAN, DualStyleGAN provides a natural way of style transfer by characterizing the content and style of a portrait with an **intrinsic style path** and a new **extrinsic style path**, respectively. The delicately designed extrinsic style path enables our model to modulate both the color and complex structural styles hierarchically to precisely pastiche the style example. Furthermore, a novel progressive fine-tuning scheme is introduced to smoothly transform the generative space of the model to the target domain, even with the above modifications on the network architecture. Experiments demonstrate the superiority of DualStyleGAN over state-of-the-art methods in high-quality portrait style transfer and flexible style

control. Code is available at <https://github.com/williamyang1991/DualStyleGAN>.

1. Introduction

Artistic portraits are popular in our daily lives and especially in industries related to comics, animations, posters, and advertising. In this paper, we focus on exemplar-based portrait style transfer, a core problem that aims to transfer the style of an exemplar artistic portrait onto a target face. Its potential application is appealing in that it allows any novice to easily transform their photograph into a stunning pastiche based on the style of their favourite artworks, which otherwise would have required highly professional skills for manual creation.

Automatic portrait style transfer based on image style transfer [23, 24, 31] and image-to-image translation [6, 19, 22] has been extensively studied. Recently, StyleGAN [17, 18], the state-of-the-art face generator, has been very promising for high-resolution artistic portrait generation via transfer learning [29]. Specifically, StyleGAN can

be effectively fine-tuned, usually only requiring hundreds of portrait images and hours of training time, to translate its generative space from the face domain to the artistic portrait domain. It shows great superiority in quality, image resolution, data requirement, and efficiency compared to image style transfer and image-to-image translation models.

The strategy above, while effective, only learns an overall translation of the distribution, incapable of performing exemplar-based style transfer. For a StyleGAN that has been transferred for generating a fixed caricature style, a laughing face will be largely mapped to its nearest one in the caricature domain, *i.e.*, a portrait with an exaggerated mouth. Users have no means of shrinking the face to pastiche their preferred artworks like in Fig. 1(c). Although StyleGAN provides inherent exemplar-based single-domain style mixing by latent swapping [1, 17], such single-domain-oriented operation is counter-intuitive and incompetent for style transfer involving a source domain and a target domain. This is because misalignment between these two domains may lead to unwanted artifacts during style mixing, especially for domain-specific structures. However, importantly, a professional pastiche, should imitate how an artist handles face structures, *e.g.*, abstraction in cartoons and deformation in caricatures.

To tackle these challenges, we propose a novel **DualStyleGAN** to realize effective modelling and control of dual styles for exemplar-based portrait style transfer. DualStyleGAN retains an **intrinsic style path** of StyleGAN to control the style of the original domain, while adding an **extrinsic style path** to model and control the style of the target extended domain, which naturally correspond to the content path and style path in the standard style transfer paradigm. Moreover, the extrinsic style path inherits the hierarchical architecture from StyleGAN to modulate structural styles in coarse-resolution layers and color styles in fine-resolution layers for flexible multi-level style manipulations.

Adding an extrinsic style path to the original StyleGAN architecture is non-trivial for our task as it risks altering the generative space and behavior of the pre-trained StyleGAN. To overcome this challenge, we present effective ways and insights to design the extrinsic style path and train DualStyleGAN. 1) *Model design*: based on the analysis on the fine-tuning behavior of StyleGAN, we propose to introduce the extrinsic style in a residual manner to the convolution layers, which can well approximate how fine-tuning affects the convolution layers of StyleGAN. We show that such design enables DualStyleGAN to effectively modulate the key structural styles. 2) *Model training*: we introduce a novel progressive fine-tuning methodology, where the extrinsic style path is first elaborately initialized so that DualStyleGAN retains the generative space of StyleGAN for seamless transfer learning. Then, we start out training DualStyleGAN with an easy style transfer task and then gradually in-

creases the task difficulty, to progressively translate its generative space to the target domain. In addition, we present a facial destylization method to provide face-portrait pairs, serving as supervision to promote the model to learn diverse styles and avoid mode collapse.

With the novel formulation above, the proposed DualStyleGAN offers high-quality and high-resolution pastiches and provides flexible and diverse control over both color styles and complicated structural styles, as shown in Fig. 1. In summary, our contributions are threefold:

- We propose a novel DualStyleGAN to characterize and control the intrinsic and extrinsic styles for exemplar-based high-resolution portrait style transfer, requiring only a few hundred style examples, which achieves superior performance over state-of-the-art methods in high-quality and diverse artistic portrait generation.
- We design a principled extrinsic style path to introduce style features from external domains via fine-tuning and to provide hierarchical style manipulation in terms of both color and structure.
- We propose a novel progressive fine-tuning scheme for robust transfer learning over networks with architecture modifications.

2. Related Work

Artistic portrait generation with StyleGAN. StyleGAN [17, 18] synthesizes high-resolution face images with hierarchical style control. Pinkney and Adler [29] fine-tuned StyleGAN on limited cartoon data, and found it promising in generating plausible cartoon faces. The original model and fine-tuned model exhibit a reasonable degree of semantic alignment [36], allowing one to toonify a real face by applying its embedded latent code in the original model to the fine-tuned model to obtain the corresponding stylized face. This framework is efficient and data-friendly, attracting further in-depth research, such as embedding acceleration [30], better choice of the latent code [34], training on extremely limited data [14, 28]. In contrast to our work, these methods only learn an overall distribution translation without exemplar-based style control. Kwong *et al.* [21] achieves style transfer by swapping fine-resolution-layer features from the exemplar style image with those from the content image under the assumption of model alignment. However, the alignment gets weakened along with unconditional fine-tuning without valid supervision, eventually leading to a failure in layer swapping. Thus, the method is mainly suitable for color transfer and not effective in controlling the vital structural styles. By comparison, our model has an explicit extrinsic style path that can be conditionally trained to characterize the structural styles. Moreover, supervision for learning diverse styles is provided via facial destylization.

Image-to-image translation. Portrait style transfer can be

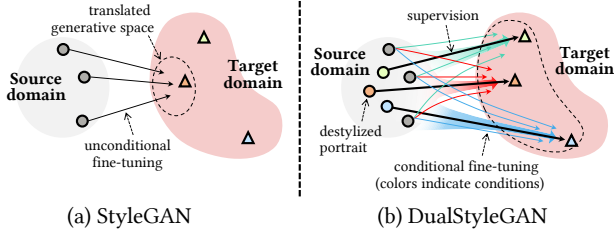


Figure 2. Compare the unconditional fine-tuning over StyleGAN and the conditional fine-tuning over DualStyleGAN.

realized by an image-to-image translation framework [27, 32, 37, 38]. The main idea is to learn a bi-directional mapping [39] between the face and artistic portrait domains. To find a correspondence between domains with a large appearance discrepancy, U-GAT-IT [19] uses attention modules to focus on key regions shared between domains. AniGAN [22] uses shared layers in discriminators to extract common features of two domains. GNR [6] learns effective content features and style features as those unchanged or altered during data augmentation, respectively. For the style of caricatures, explicit image warping is applied to imitate the distinct facial deformation [4, 33]. These strategies allow the image-to-image translation framework to stylize human faces involving drastic transformation. However, learning from scratch on complex bi-directional translations makes this framework limited to low-resolution images and requires long training time. Our method follows the fine-tuning framework of StyleGAN, which is efficient in creating high-resolution portraits and provides flexible hierarchical style control beyond the capability of above methods.

3. Portrait Style Transfer via DualStyleGAN

Our goal is to build DualStyleGAN based on a pre-trained StyleGAN, which can be transferred to a new domain and characterize the styles of both the original and the extended domains. Unconditional fine-tuning translates the StyleGAN generative space as a whole, leading to the loss of diversity of the captured styles, as illustrated in Fig. 2. Our key idea is to seek valid supervision to learn diverse styles (Sec. 3.1), and to explicitly model the two kinds of styles with two individual style paths (Sec. 3.2). We train DualStyleGAN with a principled progressive strategy for robust conditional fine-tuning (Sec. 3.3).

3.1. Facial Destylization

Facial destylization aims to recover realistic faces from artistic portraits to form anchored face-portrait pairs as supervision. Given artistic portraits of the target domain, we would like to find their reasonable counterparts in the face domain. Since the two domains might have a large appearance discrepancy, it poses us a non-trivial challenge to balance between face realism and fidelity to the portraits. To solve this problem, we propose a multi-stage destylization

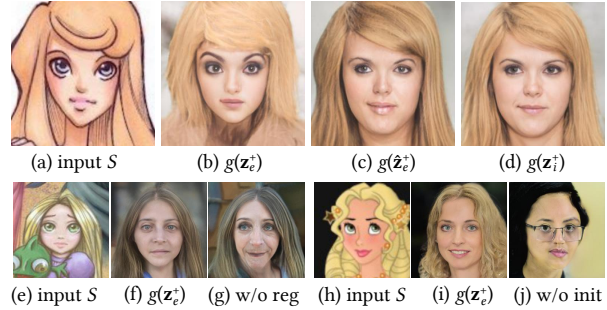


Figure 3. Illustration of facial destylization. The destylized results of (a) in each stage are sequentially shown in (b)-(d) with the exaggerated eyes gradually turning realistic. (e)-(g): Regularization prevents overfitting to the face-irrelevant green toy. (h)-(j): \mathbf{z}_e^+ serves as a good initial value to fit the complex cartoon faces.

method to gradually enhance the realism of a portrait.

Stage I: Latent initialization. The artistic portrait S is first embedded into the StyleGAN latent space by an encoder E . Here, we use a pSp encoder [30] and modify it to embed FFHQ faces [17] into \mathcal{Z}^+ space, which is more robust to face-irrelevant background details and distorted shapes than the original \mathcal{W}^+ space, as suggested in [34]. An example of the reconstructed face $g(\mathbf{z}_e^+)$ is shown in Fig. 3(b), with g the StyleGAN pre-trained on FFHQ and $\mathbf{z}_e^+ = E(S) \in \mathbb{R}^{18 \times 512}$ the latent code. Though E is trained on real faces, $E(S)$ well captures the color and the structure of portrait S . **Stage II: Latent optimization.** In [29], a face image is stylized by optimizing a latent code of g to reconstruct this image [1] and applying this code to a fine-tuned model g' . We take a reverse step to optimize the latent \mathbf{z}^+ of g' to reconstruct S with a novel regularization term, and apply the resulting $\hat{\mathbf{z}}_e^+$ to g to obtain its destylized version,

$$\hat{\mathbf{z}}_e^+ = \arg \min_{\mathbf{z}^+} \mathcal{L}_{\text{perc}}(g'(\mathbf{z}^+), S) + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(g'(\mathbf{z}^+), S) + \|\sigma(\mathbf{z}^+)\|_1, \quad (1)$$

where $\mathcal{L}_{\text{perc}}$ is the perceptual loss [15], \mathcal{L}_{ID} is the identity loss [7] to preserve the identity of the face and $\sigma(\mathbf{z}^+)$ is the standard error of 18 different 512-dimension vectors in \mathbf{z}^+ . $\lambda_{\text{ID}} = 0.1$. Different from [1], we design the regularization term to pull $\hat{\mathbf{z}}_e^+$ to the well-defined \mathcal{Z} space to avoid overfitting as in Fig. 3(f)(g), and use \mathbf{z}_e^+ rather than the mean latent code to initialize \mathbf{z}^+ before optimization, which helps accurately fit the face structures as in Fig. 3(i)(j).

Stage III: Image embedding. Finally, we embed $g(\hat{\mathbf{z}}_e^+)$ as $\mathbf{z}_i^+ = E(g(\hat{\mathbf{z}}_e^+))$, which further eliminates unreal facial details. The resulting $g(\mathbf{z}_i^+)$ has reasonable facial structures, providing valid supervision on how to deform and abstract the facial structures to imitate S .

3.2. DualStyleGAN

Figure 4 shows the network details of DualStyleGAN G . The intrinsic style path and generator network form a standard StyleGAN and are kept fixed during fine-tuning. The

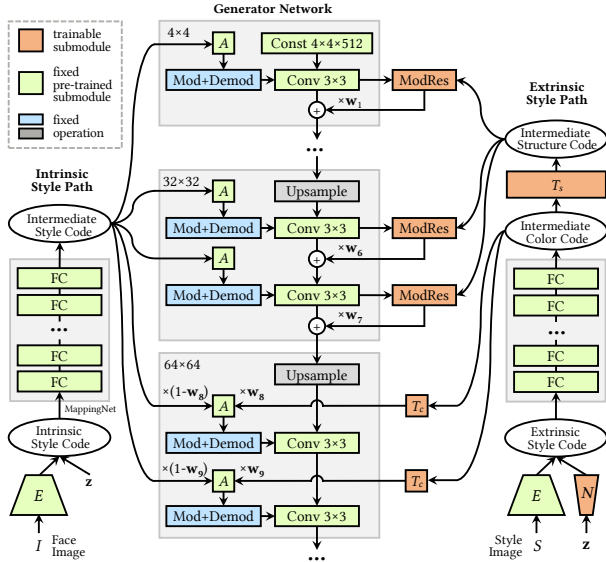


Figure 4. Network details of DualStyleGAN. For simplicity, the learned weights, biases and noises of StyleGAN are omitted.

intrinsic style path accepts intrinsic style code of unit Gaussian noise $\mathbf{z} \in \mathbb{R}^{1 \times 512}$, \mathbf{z}_i^+ of artistic portraits or \mathbf{z}^+ of real faces embedded by E . The extrinsic style path simply uses \mathbf{z}_e^+ of artistic portraits as the extrinsic style code, which captures meaningful semantic cues like hair colors and facial shapes (Fig. 3(b)). Extrinsic style codes can also be sampled via a sampling network N by mapping unit Gaussian noises to the extrinsic style distribution. Formally, given a face image I and an artistic portrait image S , exemplar-based style transfer is achieved by $G(E(I), E(S), \mathbf{w})$, where $\mathbf{w} \in \mathbb{R}^{18}$ is the weight vector for flexible style combination of two paths, and is set to $\mathbf{1}$ by default. Artistic portrait generation is realized by $G(\mathbf{z}_1, N(\mathbf{z}_2), \mathbf{w})$. When $\mathbf{w} = \mathbf{0}$, G degrades into a standard g for face generation, *i.e.*, $G(\mathbf{z}, \cdot, \mathbf{0}) \sim g(\mathbf{z})$.

StyleGAN provides a hierarchical style control, where fine-resolution and coarse-resolution layers model the low-level color style and high-level shape style, respectively, which inspires our design of the extrinsic style path.

Color control. In fine-resolution layers (8~18), the extrinsic style path takes the same strategy as StyleGAN. Specifically, \mathbf{z}_e^+ goes through a mapping network f , color transform blocks T_c and affine transform blocks A . The resulting style bias is fused with the style bias from the intrinsic style path with weight \mathbf{w} for the final AdaIN [10]. Different from g , trainable T_c composed of fully connected layers is added to characterize domain-specific colors.

Structure control. In coarse-resolution layers (1~7), we propose modulative residual blocks (ModRes) to adjust structural styles and add a structure transform block T_s to characterize domain-specific structural styles. ModRes contains a ResBlock [8] to simulate the changes of convolution layers during fine-tuning and an AdaIN block for style condition. To understand the motivation of the proposed

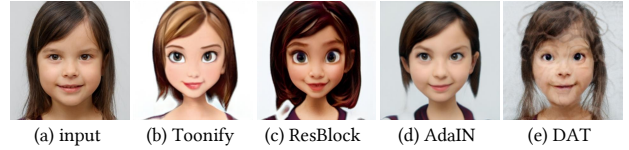


Figure 5. ResBlocks best simulate Toonify [29].

ModRes, below we provide some experimental analysis of the fine-tuning behavior over StyleGAN.

Simulating fine-tuning behavior. The success of toonification [29] relies on the semantic alignment of the models before and after fine-tuning, namely, two models have shared latent spaces [21] and closely-related convolution features. It also implies that the difference of these features is also closely-related to the original features. Moreover, among all submodules of StyleGAN, the convolution layers change the most during fine-tuning [36]. Therefore, it is possible to keep all other submodules fixed but only learn changes over the convolution features to simulate the changes of the convolution weight matrices in fine-tuning. In StyleGAN, common adjustments over deep features involve channel-wise, spatial-wise and element-wise modulations, corresponding to AdaIN [10], Diagonal Attention (DAT) [20] and ResBlock, respectively. We conduct a toy experiment and find that modulations in channel (Fig. 5(d)) or spatial (Fig. 5(e)) dimension alone are not enough to approximate the fine-tuning behavior. ResBlocks achieve the most similar results (Fig. 5(c)) to those by fine-tuning the whole StyleGAN (Fig. 5(b)). Therefore, we choose residual blocks and apply AdaIN to the convolution layers in the residual path to provide extrinsic style conditions.

Summary. Our DualStyleGAN is very simple yet effective.

- 1) **Hierarchical modelling of complex styles:** It provides hierarchical modelling on both color and complex structural styles.
- 2) **Flexible style manipulation:** It supports flexible style mixing between two domains with weight \mathbf{w} .
- 3) **Alleviate mode collapse:** Fine-tuning trains only the extrinsic style path while keeping the pre-trained StyleGAN intact, which preserves the original diverse facial features to avoid mode collapse.
- 4) **Structure preservation:** The additive property of our modulative residual block leads to a robust content loss, as we will detail in Sec. 3.3.

3.3. Progressive Fine-Tuning

We propose a progressive fine-tuning scheme to smoothly transform the generative space of DualStyleGAN towards the target domain. The scheme borrows the idea of curriculum learning [2] to gradually increase the task difficulty in three stages as illustrated in Fig. 6(a).

Stage I: Color transfer on source domain. DualStyleGAN is tasked with color transfer within the source domain in this stage. Thanks to the design of our extrinsic style path, it can be achieved purely by a specific model initialization. Specifically, the convolution filters in the modula-



(a) the goals of each stage in progressive transfer learning

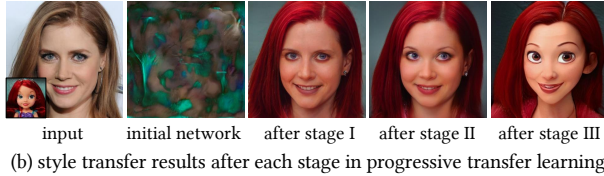


Figure 6. Illustration of progressive fine-tuning. (a) DualStyleGAN is tasked with style transfers with growing difficulties. (b) The performance of DualStyleGAN after each stage.



Figure 7. Optimize the extrinsic style code to refine color.

tive residual blocks are set to values close to 0 in order to produce negligible residual features and the fully connected layers in color transform blocks are initialized with identity matrices, meaning no changes to the input latent code. To this end, DualStyleGAN runs the standard style mixing operation of StyleGAN, where fine-resolution and coarse-resolution layers use the latent codes from the intrinsic and extrinsic style paths, respectively. As shown in Fig. 6(b), the initialized DualStyleGAN generates plausible human faces that still lie in the generative space of the pre-trained StyleGAN, allowing smooth fine-tuning in the next stage.

Stage II: Structure transfer on source domain. This stage aims to fine-tune DualStyleGAN on the source domain to fully train its extrinsic style path to capture and transfer mid-level styles. StyleGAN’s style mixing in middle layers involves small-scale style transfer like makeups, which provides DualStyleGAN with effective supervision. In stage II, we draw random latent code \mathbf{z}_1 and \mathbf{z}_2 , and would like $G(\mathbf{z}_1, \bar{\mathbf{z}}_2, \mathbf{1})$ to approximate the style mixing target $g(\mathbf{z}_l^+)$ with perceptual loss, where $\bar{\mathbf{z}}_2$ is sampled from $\{\mathbf{z}_2, E(g(\mathbf{z}_2))\}$, l is the layer where style mixing occurs and $\mathbf{z}_l^+ \in \mathcal{Z}+$ is a concatenation of l vector \mathbf{z}_1 and $(18-l)$ vector \mathbf{z}_2 . We gradually decrease l from 7 to 5 during fine-tuning with the following objective:

$$\min_G \max_D \lambda_{\text{adv}} \mathcal{L}_{\text{adv}} + \lambda_{\text{perc}} \mathcal{L}_{\text{perc}}(G(\mathbf{z}_1, \bar{\mathbf{z}}_2, \mathbf{1}), g(\mathbf{z}_l^+)), \quad (2)$$

where \mathcal{L}_{adv} is the StyleGAN adversarial loss. By decreasing

l , $g(\mathbf{z}_l^+)$ will have more structure styles from $\bar{\mathbf{z}}_2$. Thus, the extrinsic style path will learn to capture and transfer more structure styles besides colors.

Stage III: Style transfer on target domain. Finally, we fine-tune DualStyleGAN on the target domain. We would like the style codes \mathbf{z}_i^+ and \mathbf{z}_e^+ of an exemplar artistic portrait S to reconstruct S with $\mathcal{L}_{\text{perc}}(G(\mathbf{z}_i^+, \mathbf{z}_e^+, \mathbf{1}), S)$. As in the standard exemplar-based style transfer paradigm, for a random intrinsic style code \mathbf{z} , we apply style loss

$$\mathcal{L}_{\text{sty}} = \lambda_{\text{CX}} \mathcal{L}_{\text{CX}}(G(\mathbf{z}, \mathbf{z}_e^+, \mathbf{1}), S) + \lambda_{\text{FM}} \mathcal{L}_{\text{FM}}(G(\mathbf{z}, \mathbf{z}_e^+, \mathbf{1}), S),$$

where \mathcal{L}_{CX} is the contextual loss [26] and \mathcal{L}_{FM} is the feature matching loss [10], to match the style of $G(\mathbf{z}, \mathbf{z}_e^+, \mathbf{1})$ to S . For content loss, we use the identity loss [7] and L_2 regularization of the weight matrices of ModRes,

$$\mathcal{L}_{\text{con}} = \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(G(\mathbf{z}, \mathbf{z}_e^+, \mathbf{1}), g(\mathbf{z})) + \lambda_{\text{reg}} \|W\|_2. \quad (3)$$

Similar to the initialization in Stage I, regularization over weight matrices makes the residual features close to zeros, which preserves the original intrinsic facial structures and prevents overfitting. Our full objectives take the form of

$$\min_G \max_D \lambda_{\text{adv}} \mathcal{L}_{\text{adv}} + \lambda_{\text{perc}} \mathcal{L}_{\text{perc}} + \mathcal{L}_{\text{sty}} + \mathcal{L}_{\text{con}}. \quad (4)$$

3.4. Latent Optimization and Sampling

Latent optimization. It is hard to fully capture the extremely diverse styles. To solve this problem, we fix DualStyleGAN and optimize each extrinsic style code to fit its ground truth S . The optimization follows the process of embedding an image into the latent space [1] and minimizes a perceptual loss and a contextual loss in Eq. (4). As shown in Fig. 7, the colors are well refined by latent optimization.

Latent sampling. To sample random extrinsic styles, we train a sampling network N to map unit Gaussian noises to the distribution of optimized extrinsic style codes using a maximum likelihood criterion [9]. Please refer to [9] for the details. Since structures (first 7 rows of \mathbf{z}_e^+) and colors (last 11 rows of \mathbf{z}_e^+) are well disentangled in DualStyleGAN, we treat these two parts separately, *i.e.*, structure code and color code are independently sampled from N and concatenated to form the complete extrinsic style code.

4. Experiments

Datasets. Our goal is to allow users to collect portrait images of their favourite artworks for DualStyleGAN to pastiche. We would like the dataset to be limited to a few hundred images for easy collection. Therefore, we choose three datasets in popular styles of cartoon, caricature, anime. Cartoon dataset [29] has 317 images. We use 199 images from WebCaricature [11,12] and 174 images from Danbooru Portraits [3] to build the Caricature and Anime datasets, respectively. We test on the same datasets and CelebA-HQ [16,25] for extrinsic and intrinsic styles, respectively.

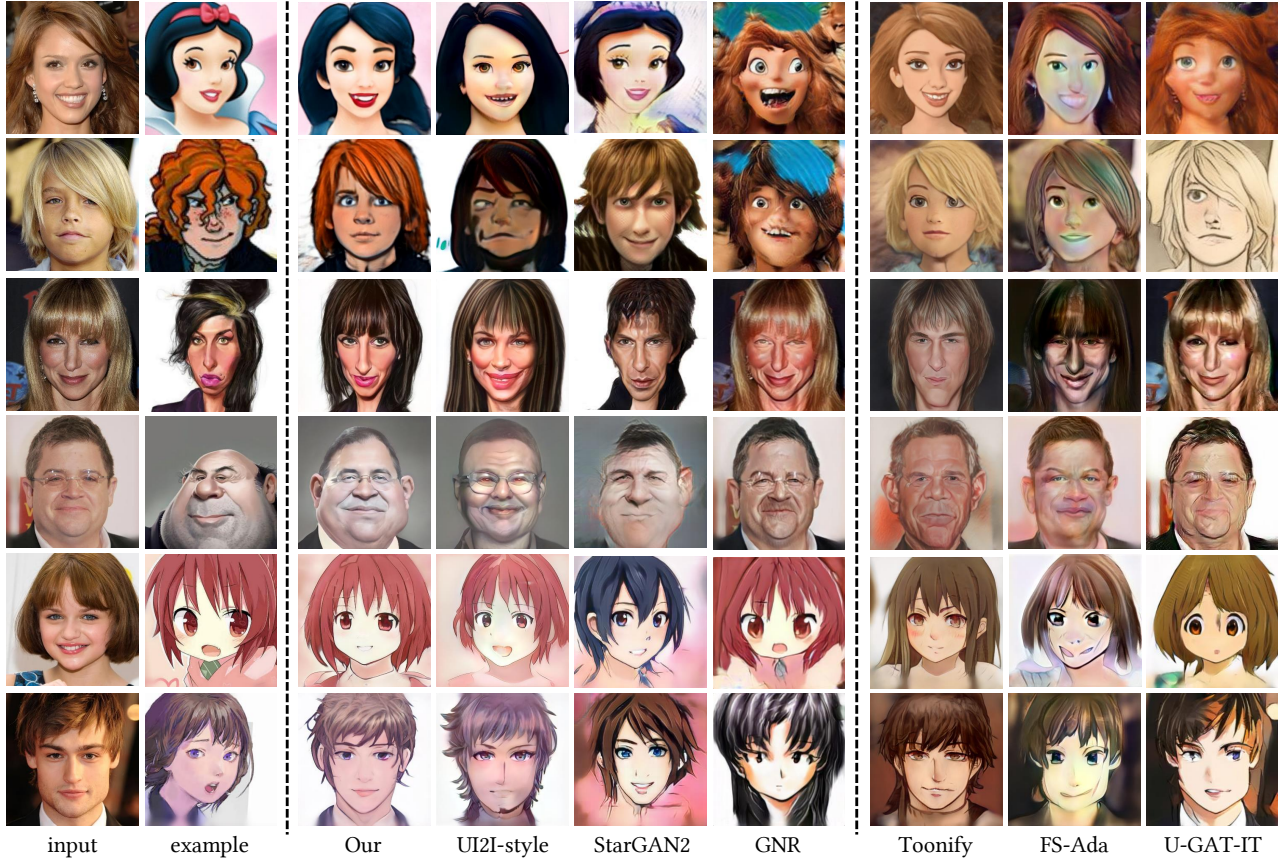


Figure 8. Visual comparison on exemplar-based portrait style transfer.

Implementation details. Our progressive fine-tuning uses eight NVIDIA Tesla V100 GPUs and a batch size of 4 per GPU. Stage II uses $\lambda_{adv} = 0.1$, $\lambda_{perc} = 0.5$, and trains on $l = 7, 6, 5$ for 300, 300, 3000 iterations, respectively, taking about 0.5 hour. Stage III sets $\lambda_{adv} = 1$, $\lambda_{perc} = 1$, $\lambda_{CX} = 0.25$, $\lambda_{FM} = 0.25$, sets $(\lambda_{ID}, \lambda_{reg})$ to $(1, 0.015)$, $(4, 0.005)$, $(1, 0.02)$ and trains for 1400, 1000, 2100 iterations on cartoon, caricature and anime, respectively. Training takes about 0.75 hour on average. Destylization (Sec. 3.1), latent optimization and training sampling network (Sec. 3.4) use one GPU and take about 5, 1, 0.13 hours, respectively. Testing takes about 0.13s per image. For simplicity, we use $[n_1 * v_1, n_2 * v_2, \dots]$ to indicate the first n_1 weights in vector \mathbf{w} are set to the value of v_1 , the next n_2 weights are set to the value of v_2 . \mathbf{w}_s and \mathbf{w}_c denote the structure weight vector (the first 7 weights of \mathbf{w}) and color weight vector (the last 11 weights), respectively. By default, we set \mathbf{w} to $\mathbf{1}$ for training and set \mathbf{w}_c to $\mathbf{1}$, \mathbf{w}_s to $\mathbf{0.75}, \mathbf{1}$ and $[4 * 0, 3 * 0.75]$ for testing for cartoon, caricature, anime, respectively.

4.1. Comparison with State-of-the-Art Methods

Figure 8 presents the qualitative comparison with six state-of-the-art methods: image-to-image-translation-based StarGAN2 [5], GNR [6], U-GAT-IT [19], and StyleGAN-

Table 1. User preference scores. Best scores are marked in bold.

Method	Cartoon	Caricature	Anime	Average
GNR [6]	0.01	0.06	0.04	0.04
StarGANv2 [5]	0.01	0.00	0.04	0.02
UI2I-style [21]	0.05	0.15	0.14	0.11
Our	0.93	0.79	0.78	0.83

based UI2I-style [21], Toonify [29], Few-Shot Adaptation (FS-Ada) [28]. The image-to-image translation and FS-Ada use 256×256 images. Other methods support 1024×1024 . Toonify, FS-Ada and U-GAT-IT learn domain-level rather than image-level styles. Thus their results are not consistent with the style examples. The severe data imbalance problem makes it hard to train valid cycle translations. Thus, StarGAN2 and GNR overfit the style images and ignore the input faces on the anime style. UI2I-style captures good color styles via layer swapping, but the model misalignment makes the structure features hard to blend, leading to failure structural style transfer, as also analyzed in Sec. 2. By comparison, DualStyleGAN transfers the best style of the exemplar style in both colors and complex structures.

To quantitatively evaluate the performance, we conduct a user study, where 27 subjects are invited to select what they consider to be the best results from the four exemplar-based

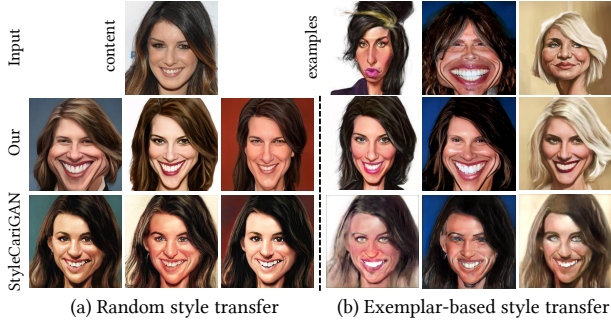


Figure 9. Comparison with StyleCariGAN [13].

style transfer methods. Each style dataset uses 10 results for evaluation. Table 1 summarizes the average preference scores, where DualStyleGAN receives the best score.

Compare with StyleCariGAN. We further compare with the advanced StyleCariGAN [13] on caricatures. StyleCariGAN combines StyleGAN and cycle translation to employ style mixing for color transfer and learn structure transfer with cycle translation. We follow it to find latent codes of the content and example images via optimization [18, 35] for its inputs. Depending on whether the latent codes are randomly sampled from the official style palette or from example images, StyleCariGAN can transfer random or exemplar-based styles on 256×256 images. As shown in Fig. 9, StyleCariGAN generates the same facial structures since its cycle translation only learns an overall structure style. By comparison, our method effectively adjusts the structure style based on the example. Moreover, our results are of higher resolution and visual quality, even if StyleCariGAN uses 6K training images, much more than ours.

4.2. Ablation Study

Paired data. Figure 10(a) compares the results with and without face-*portrait* supervision in Sec. 3.1. Without supervision, the model overfits the portraits without considering the input face structures. The supervision effectively guides the model to find the structural relationship between the face and portrait, leading to more reasonable results.

Regularization. The effect of the regularization term in the content loss (Eq. (3)) is shown in Fig. 10(b). Without the regularization term, the model overfits the hair styles of the example. Using the regularization term solves this issue. A large λ_{reg} will over-preserve the input face shape like the mouth. Therefore, we use $\lambda_{reg} = 0.005$ as a trade-off.

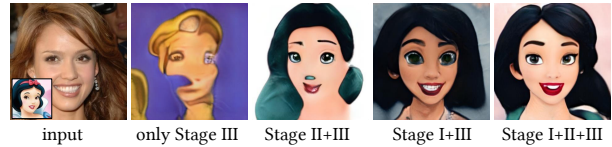
Progressive fine-tuning. As shown in Fig. 10(c), without the initialization of Stage I, the generative space of the pre-trained StyleGAN is severely altered (Fig. 6(b)), failing the transfer training completely. Without pre-training on real faces to capture face semantic features, the extrinsic style path cannot fulfill the complex task in Stage III. Only through the full progressive fine-tuning, DualStyleGAN can accurately transfer the extrinsic styles.



(a) Effect of face-*portrait* supervision: avoid overfitting



(b) Effect of the regularization term: preserve content



(c) Effect of progressive transfer learning: robustness

Figure 10. Ablation study.

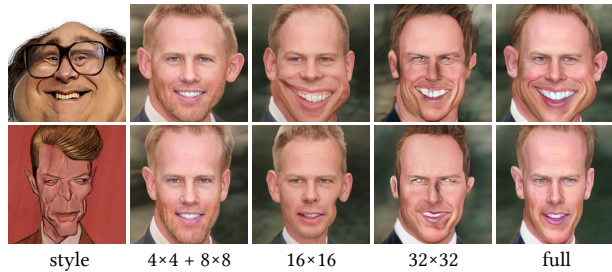


Figure 11. The proposed extrinsic style path learns semantically hierarchical structure modulations.

Effect of different layers. To study how each layer of the extrinsic style path affects the facial features, each time we activate a subset of layers (for example, $w = [3 * 0, 2 * 1, 13 * 0]$ only activates two 16×16 layers) and compare their results in Fig. 11. Since the AdaIN-based color modulation has been well studied in StyleGAN, we only focus on the structure modulation in coarse-resolution layers. It can be seen that the initial layers adjust the overall face shapes, 16×16 layers exaggerate facial components like mouths, and 32×32 layers focus on local shapes like wrinkles.

4.3. Further Analysis

Color and structure preservation. Users may want to keep the color of the original photo as Toonify [29]. We provide two ways of color preservation. The first is just to deactivate color-related layers in the extrinsic style path by setting $w_c = 0$ as in Fig. 12(c). Another way is to replace the extrinsic style codes with intrinsic style codes in the last 11 layers. Compared to the first way, the intrinsic latent codes additionally go through color transform blocks, making the final color more aligned with the target domain as in Fig. 12(d). Finally, structure preservation can be easily achieved by setting $w_s < 1$. Figure 12(e) presents exam-



Figure 12. Preservation of color and structure from the photo.

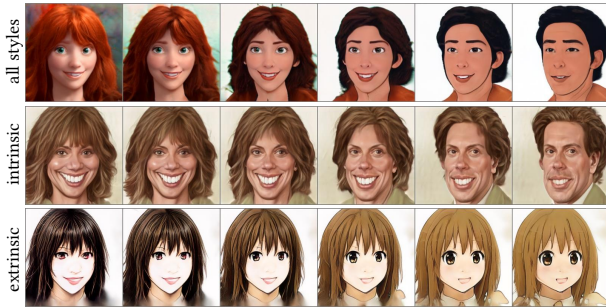


Figure 13. Intrinsic and extrinsic style blending.

ples of mild style transfer with $w_s = 0.5$.

Style blending. In Fig. 13, we fuse styles by interpolating two intrinsic and/or extrinsic style codes. The smooth transition implies a reasonable coverage of the style manifold.

Performance on other styles. We further collect datasets in styles of Pixar, Comic and Slam Dunk from the Internet, with 122, 101 and 120 images, respectively. Our method achieves good performance on these styles as in Fig. 14.

Performance on unseen style. Given unseen styles beyond training data, our method still transfers reasonable but less consistent styles (Fig. 15(c)). By destylizing the unseen image to obtain a fixed intrinsic style code and optimizing the extrinsic style code as in Sec 3.4, better styles are learned (Fig. 15(d)). However, it introduces some artifacts. We leave robust unseen style extensions to future work.

4.4. Limitations

In Fig. 16, we show three typical failure cases of DualStyleGAN. First, while the face features are well captured, details in non-facial regions like the hat and background textures are lost in our result. Second, anime faces often have very abstract noses. If we retain the color of the photo, the nose becomes evident but unnatural for anime style. Third, our method still suffers from data bias problem. The anime dataset has strong bias towards straight hairs and bangs, making our method fail to handle curly hairs without bangs. Meanwhile, uncommon styles like the extremely large eyes cannot be well imitated. As a result, applying our method to tasks with severe data imbalance problem might lead to unsatisfactory results on under-representated data.

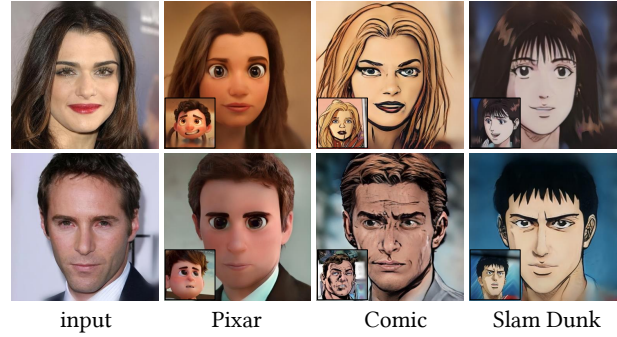


Figure 14. Performance on Pixar, Comic and Slam Dunk styles.



Figure 15. Performance on unseen style.

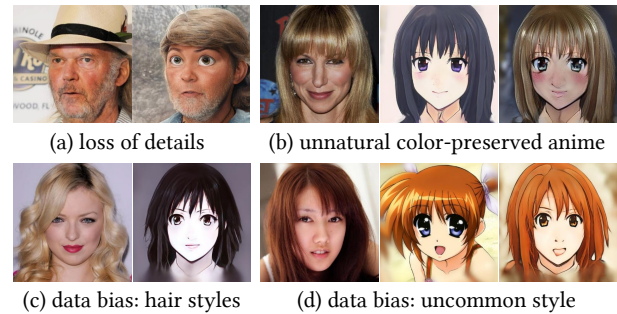


Figure 16. Limitations of DualStyleGAN

5. Conclusion and Future Work

In this paper, we extend StyleGAN to accept style condition from new domains while preserving its style control in the original domain. This results in an interesting application of high-resolution exemplar-based portrait style transfer with a friendly data requirement. DualStyleGAN, with an additional style path to StyleGAN, can effectively model and modulate the intrinsic and extrinsic styles for flexible and diverse artistic portrait generation. We show that valid transfer learning on DualStyleGAN can be achieved with especial architecture design and progressive training strategy. We believe our idea of model extension in terms of both architecture and data can be potentially applied to other tasks such as more general image-to-image translation and knowledge distillation. In future work, we would like to explore the recommendation of the suitable style image and its weight vector w for the input photo for easy use and to alleviate the data bias problem via data augmentation.

Acknowledgments. This study is supported under the RIE2020 Industry Alignment Fund – Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s).

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proc. Int'l Conf. Computer Vision*, pages 4432–4441, 2019. 2, 3, 5
- [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proc. IEEE Int'l Conf. Machine Learning*, pages 41–48, 2009. 4
- [3] Gwern Branwen, Anonymous, and Danbooru Community. Danbooru2019 portraits: A large-scale anime head illustration dataset. <https://www.gwern.net/Crops#danbooru2019-portraits>, March 2019. 5
- [4] Kaidi Cao, Jing Liao, and Lu Yuan. Carigans: unpaired photo-to-caricature translation. *ACM Transactions on Graphics*, 37(6):1–14, 2018. 3
- [5] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 8188–8197, 2020. 6
- [6] Min Jin Chong and David Forsyth. GANs N' Roses: Stable, controllable, diverse image to image translation. *arXiv preprint arXiv:2106.06561*, 2021. 1, 3, 6
- [7] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 4690–4699, 2019. 3, 5
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 770–778, 2016. 4
- [9] Yedid Hoshen, Ke Li, and Jitendra Malik. Non-adversarial image synthesis with generative latent nearest neighbors. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 5811–5819, 2019. 5
- [10] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proc. Int'l Conf. Computer Vision*, pages 1510–1519, 2017. 4, 5
- [11] Jing Huo, Yang Gao, Yinghuan Shi, and Hujun Yin. Variation robust cross-modal metric learning for caricature recognition. In *Proc. Thematic Workshops of ACM Int'l Conf. Multimedia*, pages 340–348, 2017. 5
- [12] Jing Huo, Wenbin Li, Yinghuan Shi, Yang Gao, and Hujun Yin. Webcaricature: a benchmark for caricature recognition. In *Proc. British Machine Vision Conference*, 2018. 5
- [13] Wonjong Jang, Gwangjin Ju, Yuchool Jung, Jiaolong Yang, Xin Tong, and Seungyong Lee. StyleCariGAN: caricature generation via stylegan feature map modulation. *ACM Transactions on Graphics*, 40(4):1–16, 2021. 7
- [14] Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy. Deceive D: Adaptive pseudo augmentation for gan training with limited data. In *Advances in Neural Information Processing Systems*, 2021. 2
- [15] Justin Johnson, Alexandre Alahi, and Fei Fei Li. Perceptual losses for real-time style transfer and super-resolution. In *Proc. European Conf. Computer Vision*, pages 694–711. Springer, 2016. 3
- [16] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *Proc. Int'l Conf. Learning Representations*, 2018. 5
- [17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1, 2, 3
- [18] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 1, 2, 7
- [19] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwang Hee Lee. U-GAT-IT: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In *Proc. Int'l Conf. Learning Representations*, 2019. 1, 3, 6
- [20] Gihyun Kwon and Jong Chul Ye. Diagonal attention and style-based gan for content-style disentanglement in image generation and translation. In *Proc. Int'l Conf. Computer Vision*, 2021. 4
- [21] Sam Kwong, Jialu Huang, and Jing Liao. Unsupervised image-to-image translation via pre-trained stylegan2 network. *IEEE Transactions on Multimedia*, 2021. 2, 4, 6
- [22] Bing Li, Yuanlue Zhu, Yitong Wang, Chia-Wen Lin, Bernard Ghanem, and Linlin Shen. AniGAN: Style-guided generative adversarial networks for unsupervised anime face generation. *IEEE Transactions on Multimedia*, 2021. 1, 3
- [23] Chuan Li and Michael Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 2479–2486, 2016. 1
- [24] Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. Visual attribute transfer through deep image analogy. *ACM Transactions on Graphics*, 36(4):120, 2017. 1
- [25] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proc. Int'l Conf. Computer Vision*, pages 3730–3738, 2015. 5
- [26] Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor. The contextual loss for image transformation with non-aligned data. In *Proc. European Conf. Computer Vision*, pages 768–783, 2018. 5
- [27] Ori Nizan and Ayellet Tal. Breaking the cycle-colleagues are all you need. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 7860–7869, 2020. 3
- [28] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. Few-shot image generation via cross-domain correspondence. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 10743–10752, 2021. 2, 6
- [29] Justin NM Pinkney and Doron Adler. Resolution dependent gan interpolation for controllable image synthesis between domains. *arXiv preprint arXiv:2010.05334*, 2020. 1, 2, 3, 4, 5, 6, 7
- [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 *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 2021. 2, 3
- [31] Ahmed Selim, Mohamed Elgharib, and Linda Doyle. Painting style transfer for head portraits using convolutional neural networks. *ACM Transactions on Graphics*, 35(4):1–18, 2016. 1
- [32] Xuning Shao and Weidong Zhang. SPatchGAN: A statistical feature based discriminator for unsupervised image-to-image translation. In *Proc. Int'l Conf. Computer Vision*, 2021. 3
- [33] Yichun Shi, Debayan Deb, and Anil K Jain. WarpGAN: Automatic caricature generation. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 10762–10771, 2019. 3
- [34] 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 Transactions on Graphics*, 40(4):1–13, 2021. 2, 3
- [35] Ayush Tewari, Mohamed Elgharib, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. PIE: Portrait image embedding for semantic control. *ACM Transactions on Graphics*, 39(6):1–14, 2020. 7
- [36] Zongze Wu, Yotam Nitzan, Eli Shechtman, and Dani Lischinski. StyleAlign: Analysis and applications of aligned stylegan models. *arXiv preprint arXiv:2110.11323*, 2021. 2, 4
- [37] Shaoan Xie, Mingming Gong, Yanwu Xu, and Kun Zhang. Unaligned image-to-image translation by learning to reweight. In *Proc. Int'l Conf. Computer Vision*, pages 14174–14184, 2021. 3
- [38] Yihao Zhao, Ruihai Wu, and Hao Dong. Unpaired image-to-image translation using adversarial consistency loss. In *Proc. European Conf. Computer Vision*, pages 800–815. Springer, 2020. 3
- [39] Jun Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proc. Int'l Conf. Computer Vision*, pages 2242–2251, 2017. 3