

# Neural Hair Rendering

Menglei Chai, Jian Ren, and Sergey Tulyakov

Snap Inc.

**Abstract.** In this paper, we propose a generic neural-based hair rendering pipeline that can synthesize photo-realistic images from virtual 3D hair models. Unlike existing supervised translation methods that require model-level similarity to preserve consistent structure representation for both real images and fake renderings, our method adopts an unsupervised solution to work on arbitrary hair models. The key component of our method is a shared latent space to encode appearance-invariant structure information of both domains, which generates realistic renderings conditioned by extra appearance inputs. This is achieved by domain-specific pre-disentangled structure representation, partially shared domain encoder layers, and a structure discriminator. We also propose a simple yet effective temporal conditioning method to enforce consistency for video sequence generation. We demonstrate the superiority of our method by testing it on large amount of portraits, and comparing with alternative baselines and state-of-the-art unsupervised image translation methods.

**Keywords:** Neural rendering, unsupervised image translation

## 1 Introduction

Hair is a critical component of human subjects. Rendering virtual 3D hair models into realistic images has been long studied in computer graphics, due to the extremely complicated geometry and material of human hair. Traditional graphical rendering pipeline simulates every aspect of natural hair appearance, including surface shading, light scattering, semi-transparent occlusions and soft shadowing by leveraging physics-based shading models of hair fibers, global illumination rendering algorithms to capture mutual interactions between fibers and the environment, and artistically designed material parameters. Given the extreme complexity of the geometry and associated lighting effects of hair, such a direct approximation of physical hair appearance requires a highly detailed 3D model, carefully tuned material parameters, and huge amount of rendering computation, which are often too costly and unaffordable for interactive application scenarios that require efficient feedback and user-friendly interactions, such as games and real-time photo editing software.

With the recent advances in generative adversarial networks, it becomes natural to formulate hair rendering as a special case of the conditional image generation problem, with the hair structure controlled by the 3D model, while realistic appearance is synthesized by a neural network. In the context of image-to-image

translation, one of the major challenges is how to bridge both the source and target domains for proper translation. Most existing hair generation methods fall into the supervised category, which demands enough training image pairs to provide direct supervision. For example, sketch-based hair generation methods [25,19,38] construct training pairs by synthesizing user sketches from real images. While a number of such methods were introduced, rendering 3D hair models with the help of neural networks did not receive similar treatment. The existing work on the topic [54] requires real and fake domains considerably overlap, such that the common structure is present in both domains. This is achieved at the cost of a very complicated strand-level high quality model, and allows for extracting edge and orientation maps from rendered hair strands, which serve as the common representations of hair structures between real photos and fake models. However, preparing such a high-quality strand-level hair model is itself an expensive and non-trivial problem even for a professional artist, which significantly restricts the application scope of this method.

In this paper, we propose a generic neural-network-based hair rendering pipeline that provides efficient and realistic rendering of a generic low-quality 3D hair model borrowing the material features extracted from an arbitrary reference hair image. Instead of using a complicated strand-level models to match real-world hairs like [54], we allow users to use any type of hair model requiring only the isotropic structure of hair strands be properly represented. Particularly, we adopt sparse polygon strip meshes which are much more widely used in interactive applications [53]. Given the dramatic difference between such a coarse geometry and real hair, we are not able to design common structure representations at the model level. Therefore, supervised image translation methods will be infeasible due to the lack of paired data.

To bridge the domain of real hair images and low-quality virtual hair models in an unsupervised manner, we propose to construct a shared latent space between both real and fake domains, which encodes a common structural representation even if inputs from different domains are totally different, and render the realistic hair image from this latent space with the appearance conditioned by an extra input. This is achieved by: 1) different domain structure encodings are used as the network inputs, to pre-disentangle geometric structure and chromatic appearance for both real hair images and 3D models; 2) a UNIT [30]-like architecture is adopted to enable common latent space by partially sharing encoder weights between two auto-encoder branches that are trained with in-domain supervision; 3) a structure discriminator is introduced to further match the distribution of the encoded structure features; 4) supervised reconstruction is enforced on both branches to guarantee all necessary structure information is kept in the shared feature space. In addition, to enable temporally-smooth animation rendering, we introduce a simple yet effective temporal condition method with single image training data only, utilizing the exact hair model motion fields. We demonstrate the effectiveness of the pipeline and each key component by extensively testing on a large amount of diverse human portraits and various hair models. We also compare our method with general unsupervised image transla-

tion methods, and show that due to the limited sampling ability on the synthetic hair domain, all existing methods fail to produce convincing results.

## 2 Related Work

Conditional neural hair rendering belongs to a wide range of problems tackling portrait manipulation and editing. A number of methods in the literature address this cross-domain generation problem such as paired and unpaired image-to-image translation and style transfer.

**Image-to-image translation** aims at converting images from one domain to another while keeping the structure of the source image unchanged. The literature contains a number of various methods performing this task in a variety of settings. Paired image-to-image translation methods [18,52] operate when pairs of images in the source and the target domains are available. Such methods, for example, translate from semantic labels to scenes [52,37,2], from edges to objects [43], and perform image super-resolution [24,20]. However, paired data are not always available in many practical applications. Unsupervised image-to-image translation tackles a setting in which paired data is not available, while sampling from two domains is possible [31,47,58,6,44,30,17]. Clearly, unpaired image-to-image translation is an ill-posed problem for there are numerous ways an image can be transformed to a different domain. Hence, recently proposed methods introduce constraints to limit the number of possible transformations. Some studies enforce certain domain properties [1,44], while other concurrently introduced works apply cycle-consistency to transform images between different domains, such as horse to zebra, day-to-night [57,58,22]. Our work differs from existing studies that we focus on a specific challenging problem, which is the realistic hair generation, where we want to translate manually designed hair models from the domain of rendered images to the domain of real hair. For the purpose of controllable hair generation, we leverage rendered hair structure and arbitrary hair appearance to synthesize diverse realistic hair styles. Further difference of our work compared to the image-to-image translation papers is unbalanced data. The domain of images containing real hair is far more diverse than that of rendered hair, making it challenging for classical image-to-image translation works to address the problem.

**Neural style transfer and manipulation** is related to image-to-image translation in a way that image style is changed while content is maintained [3,10,16,27,29,28,51,13]. Style in this case is represented by unique style of an artist [10,51] or is copied from an example image provided by the user. Our work follows the research idea from example-guide style transfer that hair style is obtained from reference real image. However, instead of changing style of a whole image, our aim is to keep the appearance of human face and background unchanged, while having full control over the hair region. Therefore, instead of following exiting works that inject style features into image generation networks directly [16,37], we propose a new architecture that combines only hair appearance features and latent features that encodes image content and adapted hair

structure for image generation. This way we can achieve the goal that only the style of the hair region is manipulated according to provided exemplar image.

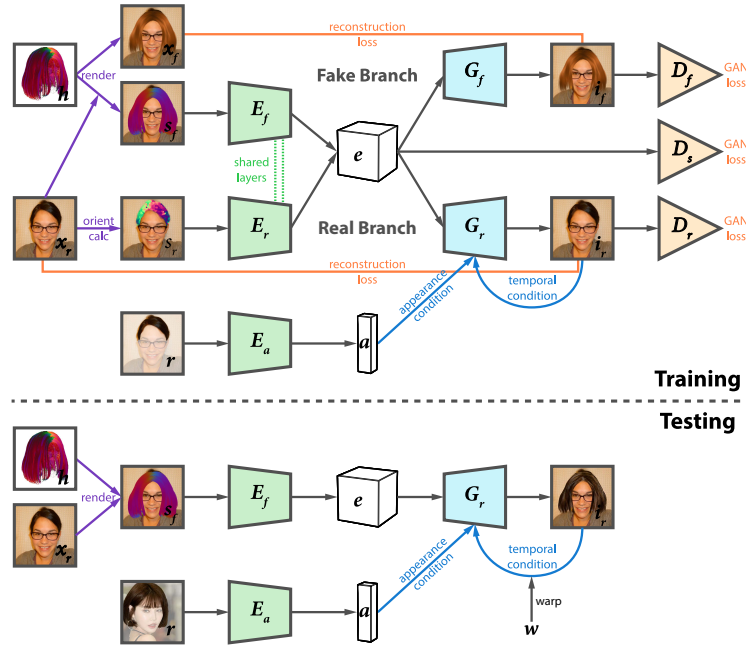
**Domain Adaptation** addresses the domain-shift problem that is widely exists between source and target domains [42]. Various feature based methods have been proposed to tackle the problem [23,11,12,7,50]. Recent works on adversarial learning for the embedded feature alignment between source and target domains achieves better results than previous studies [8,9,32,48,15,49]. Efforts using domain adaptation for both classification and pixel-level prediction tasks have gained significantly progress [1,4,48]. In this work, we follow the challenging setting of unsupervised domain adaptation that there is no corresponding annotation between source and target domains. We aim at learning an embedding space that only contains hair structure information for both rendered and real domain. Considering the domain gap, instead of using original images as input, we use rendered and real structure map as inputs to the encoders, which contain both domain specific layers and shared layers, to obtain latent features. The adaptation is achieved by adversarial training and image reconstruction.

**Hair Rendering and Generation** share a similar goal with our paper, which is synthesizing photo-realistic hair images. Traditional graphical hair rendering methods focus on improving rendering quality and performance by either more accurately modeling the special hair material and lighting behaviours [33,34,5,56], or approximating certain aspects of rendering pipeline to reduce the computation complexity [59,35,41,39,55]. However, the extremely huge computation cost for realistic hair rendering usually prohibits them to be directly applied in real-time applications. Utilizing the latest advances in GANs, recent works [25,19,38] achieved impressive progress on conditioned hair image generation as supervised image-to-image translation. A GAN-based hair rendering method [54] proposes to perform conditioned 3D hair rendering by starting with a common structure representation and progressively enforce various conditions. However, it requires the hair model to be able to generate consistent representation (strand orientation map) with real images, which is challenging for low-quality mesh-based models, and cannot achieve temporally smooth results.

Despite recent progress on conditional hair generation and image-to-image translation, the problem of realistic spatio-temporal rendering of low-quality 3D hair models remains largely unaddressed. In this paper we provide the necessary treatment of the problem achieving photo-realistic renderings as well as reach temporal stability of the rendered hair.

### 3 Approach

**Problem Formulation.** Let  $\mathbf{h}$  be the target 3D hair model, with camera parameters  $\mathbf{c}$  and hair material parameters  $\mathbf{m}$ , we formulate the *traditional graphic rendering pipeline* as  $R_t(\mathbf{h}, \mathbf{m}, \mathbf{c})$ . Likewise, our *neural network-based rendering pipeline* is defined as  $R_n(\mathbf{h}, \mathbf{r}, \mathbf{c})$ , with a low-quality hair model  $\mathbf{h}$  and material features extracted from an arbitrary reference hair image  $\mathbf{r}$ .



**Fig. 1. The overall pipeline of our neural hair rendering framework.** We use two branches to encode hair structure features, where one for the real domain and another for the fake domain. A domain discriminator is applied to the latent space to find domain invariant features. We also use two decoders to reconstruct images for two domains. The decoder in the real domain is different from the one in the fake domain, for it is conditioned on a reference image. Additionally, to generate consistent videos, we apply a temporal condition on the real branch. During inference, we use the encoder in the fake branch to get hair structure features from a 3D hair model and use the generator in the real branch to synthesized an appearance conditioned image.

### 3.1 Overview of Network Architecture

The overall system pipeline is shown in Fig.1, which consists of two parallel branches for both domains of real photo (i.e., *real*) and synthetic renderings (i.e., *fake*), respectively.

On the encoding side, the *structure adaptation subnetwork*, which includes a real encoder  $E_r$  and a fake encoder  $E_f$ , achieves cross-domain structure embedding  $e$ . Similar to UNIT[30], we share the weights of the last few ResNet layers in  $E_r$  and  $E_f$  to extract consistent structural representation from two domains. In addition, a structure discriminator  $D_s$  is introduced to match the high-level feature distributions between two domains to enforce the shared latent space further to be domain invariant.

On the decoding side, the *appearance rendering subnetwork*, which consists of  $G_r$  and  $G_f$  for the real and fake domain respectively, is attached after the shared latent space  $e$  to reconstruct the images in the corresponding domain.

Each decoder owns its exclusive domain discriminator  $D_r$  and  $D_f$  to ensure the reconstruction matches the domain distribution, besides the reconstruction losses. The hair appearance is conditioned in an asymmetric way that  $G_r$  accepts extra condition of material features extracted from a reference image  $\mathbf{r}$  by using material encoder  $E_m$ , while the unconditional decoder  $G_f$  is asked to memorize the appearance, which is made on purpose for training data generation (Sec.4.1).

At the *training stage*, all these networks are jointly trained using two sets of image pairs  $(\mathbf{s}, \mathbf{x})$  for both real and fake domains, where  $\mathbf{s}$  represents a domain-specific structure representation of the corresponding hair image  $\mathbf{x}$  in this domain. Both real and fake branches try to reconstruct the image  $G(E(\mathbf{x}))$  from its paired structure image  $\mathbf{s}$  independently through their own encoder-decoder networks, while the shared structural features are enforced to match each other consistently by the structure discriminator  $D_s$ . We set the appearance reference  $\mathbf{r} = \mathbf{x}$  in the real branch to fully reconstruct  $\mathbf{x}$  in a paired manner.

At the *inference stage*, only the fake branch encoder  $E_f$  and the real branch decoder  $G_r$  are activated.  $G_r$  generates the final realistic rendering using structural features encoded by  $E_f$  on the hair model. The final rendering equation  $R_n$  can be formulated as:

$$R_n(\mathbf{h}, \mathbf{r}, \mathbf{c}) = G_r(E_f(S_f(\mathbf{h}, \mathbf{c})), E_m(\mathbf{r})), \quad (1)$$

where the function  $S_f(\mathbf{h}, \mathbf{c})$  renders the structure encoded image  $\mathbf{s}_f$  of the model  $\mathbf{h}$  in camera setting  $\mathbf{c}$ .

### 3.2 Structure Adaptation

The goal of the structure adaptation subnetwork, formed by the encoding parts of both branches, is to encode cross-domain structural features to support final rendering. Since the inputs to both encoders are manually disentangled structure representation (Sec.4.1), the encoded features  $E(\mathbf{s})$  only contain structural information of the target hair. Moreover, as the appearance information is either conditioned by extra decoder input in a way that non spatial-varying structural information is leaked (the real branch) or simple enough to be memorized by the decoder (the fake branch) (Sec.3.3), the encoded features should also include all the structural information necessary to reconstruct  $\mathbf{x}$ .

$E_r$  and  $E_f$  share a similar network structure: five downsampling convolution layers followed by six ResBlks. The last two ResBlks are weight-sharing to enforce the shared latent space.  $D_s$  follows PatchGAN[18] to distinguish between the latent feature maps from both domains:

$$\mathcal{L}_{D_s} = \mathbb{E}_{\mathbf{s}_r} [\log(D_s(E_r(\mathbf{s}_r)))] + \mathbb{E}_{\mathbf{s}_f} [\log(1 - D_s(E_f(\mathbf{s}_f)))] \quad (2)$$

### 3.3 Appearance Rendering

The hair appearance rendering subnetwork decodes the shared cross-domain hair features into the real domain images. The decoders  $G_r$  and  $G_f$  have different

network structures and do not share weights since the neural hair rendering is a unidirectional translation that aims to map the rendered 3D model in the fake domain to real images in the real domain. Therefore,  $G_f$  is required to make sure the latent features  $\mathbf{e}$  encode all necessary information from the input 3D model, instead of learning to render various appearance. On the other hand,  $G_r$  is designed in a way to accept arbitrary inputs for realistic image generation.

Specifically, the unconditional decoder  $G_f$  starts with two ResBlks, and then five consecutive upsampling transposed convolutional layers followed by one final convolutional layer.  $G_r$  adopts a similar structure as  $G_f$ , with each transposed convolutional layer replaced with a SPADE[37] ResBlk to use appearance feature maps  $\mathbf{a}_{\mathbf{r},\mathbf{s}_r}$  at different scales to condition the generation. Assuming the binary hair mask of the reference and the target images are  $\mathbf{m}_{\mathbf{r}}$  and  $\mathbf{m}_{\mathbf{s}}$ , the appearance encoder  $E_m$  extracts the appearance feature vector on  $\mathbf{r} \times \mathbf{m}_{\mathbf{r}}$  with five downsampling convolutional layers and an average pooling. This feature vector  $E_m(\mathbf{r})$  is then used to construct the feature map  $\mathbf{a}_{\mathbf{r},\mathbf{s}_r}$  by duplicating it spatially in the target hair mask  $\mathbf{m}_{\mathbf{s}}$  as follows:

$$\mathbf{a}_{\mathbf{r},\mathbf{s}_r}(p) = \begin{cases} E_m(\mathbf{r}), & \text{if } \mathbf{m}_{\mathbf{s}_r}(p) = 1, \\ 0, & \text{if } \mathbf{m}_{\mathbf{s}_r}(p) = 0. \end{cases} \quad (3)$$

To make sure the reconstructed real image  $G_r(E_r(\mathbf{s}_r), \mathbf{a}_{\mathbf{r},\mathbf{s}_r})$  and the reconstructed fake image  $G_f(E_f(\mathbf{s}_f))$  belong to their respective distributions, we apply two domain specific discriminator  $D_r$  and  $D_f$  for the real and fake domain respectively. The adversarial losses write as:

$$\mathcal{L}_{D_r} = \mathbb{E}_{\mathbf{x}_r}[\log(D_r(\mathbf{x}_r))] + \mathbb{E}_{\mathbf{s}_r, \mathbf{r}}[\log(1 - D_r(G_r(E_r(\mathbf{s}_r), \mathbf{a}_{\mathbf{r},\mathbf{s}_r})))] \quad (4)$$

$$\mathcal{L}_{D_f} = \mathbb{E}_{\mathbf{x}_f}[\log(D_f(\mathbf{x}_f))] + \mathbb{E}_{\mathbf{s}_f}[\log(1 - D_f(G_f(E_f(\mathbf{s}_f))))] \quad (5)$$

We also adopt perceptual losses to measure high-level feature distance utilizing the paired data:

$$\mathcal{L}_p = \sum_{l=0}^L \|\Psi_l(G_r(E_r(\mathbf{s}_r), \mathbf{a}_{\mathbf{r},\mathbf{s}_r})) - \Psi_l(\mathbf{x}_r)\|_1 + \|\Psi_l(G_f(E_f(\mathbf{s}_f))) - \Psi_l(\mathbf{x}_f)\|_1, \quad (6)$$

where  $\Psi_l(\mathbf{i})$  computes the activation feature map of input image  $\mathbf{i}$  at the  $l$ th selected layer of VGG-19[45] pre-trained on ImageNet[40].

Finally, we have the overall training objective as:

$$\min_{E,G} \max_D (\lambda_s \mathcal{L}_{D_s} + \lambda_g (\mathcal{L}_{D_r} + \mathcal{L}_{D_f}) + \lambda_p \mathcal{L}_p). \quad (7)$$

### 3.4 Temporal Conditioning

The aforementioned rendering network is able to generate plausible single-frame results. However, despite the hair structure is controlled by smoothly-varying

inputs of  $\mathbf{s}_f$  with the appearance conditioned by a fixed feature map  $\mathbf{a}_{r,s_r}$ , the spatially-varying appearance details are still generated in a somewhat arbitrary manner which tends to flicker in time (Fig.5). Fortunately, with the availability of the 3D model, we can calculate the exact hair motion flow  $\mathbf{w}^t$  for each pair of frames  $t-1$  and  $t$ , which can be used to warp image  $\mathbf{i}$  from  $t-1$  to  $t$  as  $W(\mathbf{i}, \mathbf{w}^t)$ . We utilize this dense correspondences to enforce temporal smoothness.

Let  $\mathbf{I} = \{\mathbf{i}^0, \mathbf{i}^1, \dots, \mathbf{i}^T\}$  be the generated result sequence, we achieve this temporal conditioning by simply using the warped result of the previous frame  $W(\mathbf{i}^{t-1}, \mathbf{w}^t)$  as an additional condition, stacked with the appearance feature map  $\mathbf{a}_{r,s_r}$ , to the real branch decoder  $G_r$  when generating the current frame  $\mathbf{i}^t$ .

We make the network temporally consistent by changing the real branch decoder only. Specifically, we temporally finetune it. During temporal training, we fix all other networks and use the same objective as Eq.7, but randomly (50% of chance) concatenating  $\mathbf{x}_r$  into the condition inputs to the SPADE ResBlks of  $G_r^t$ . The generation pipeline of the real branch now becomes  $G_r^t(E_r(\mathbf{s}_r), \mathbf{a}_{r,s_r}, \mathbf{x}_r)$ , so that the network learns to preserve the temporal consistency if the previous frame is inputted as the temporal condition, or generate randomly from scratch if the temporal condition is set to zero.

Finally, we have the rendering equation for sequential generation:

$$\mathbf{i}^t = R_n(\mathbf{h}, \mathbf{r}, \mathbf{c}^t) = \begin{cases} G_r(E_f(\mathbf{s}_f^t), \mathbf{a}_{r,s_f^t}), & \text{if } t = 0, \\ G_r^t(E_f(\mathbf{s}_f^t), \mathbf{a}_{r,s_f^t}, W(\mathbf{i}^{t-1}, \mathbf{w}^t)). & \text{if } t > 0, \end{cases} \quad (8)$$

$$\mathbf{s}_f^t = S_f(\mathbf{h}, \mathbf{c}^t).$$

## 4 Experiments

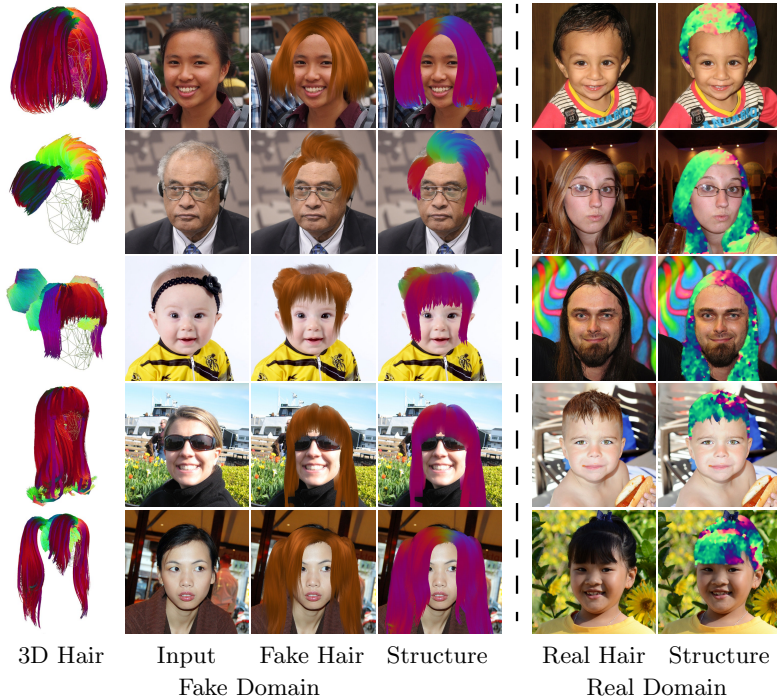
In this section, we show the experimental results of our proposed neural hair rendering and demonstrate its superiority over existing state-of-the-art works.

### 4.1 Data Preparation

To train the proposed framework, we generate a dataset that includes image pairs  $(\mathbf{s}, \mathbf{x})$  for both real and fake domains. In each domain,  $\mathbf{s} \rightarrow \mathbf{x}$  indicates the mapping from structure to image, where  $\mathbf{s}$  encodes only the structure information, and  $\mathbf{x}$  is the corresponded image that conforms to the structure condition.

**Real Domain.** We adopt the widely used FFHQ[21] portrait dataset to generate the training pairs for the real branch, given it contains diverse hairstyles on shapes and appearances. To prepare real data pairs, we use original portrait photos from FFHQ as  $\mathbf{x}_r$ , and generate  $\mathbf{s}_r$  to encode only structure information from hair. However, obtaining  $\mathbf{s}_r$  is a non-trivial process since hair image also contains material information, besides structural knowledge. To fully disentangle structure and material, and construct a universal structural representation  $\mathbf{s}$  of all real hair, we apply a dense pixel-level orientation map in the hair region, which is formulated as  $\mathbf{s}_r = S_r(\mathbf{x}_r)$ , calculated with oriented filter kernels [36].





**Fig. 2. Training data preparation.** We prepare two groups of training data. For the fake domain (a), we use hair model and input image to generate both fake rendering and model structure map. For the real domain (b), we generate image structure map for each image.

Thus, we can obtain  $\mathbf{s}_r$  that only consists of local hair strand flow structures. Example generated pairs are presented in Fig.2b.

For the purpose of training and validation, we randomly select 65,000 images from FFHQ as training, and use the remaining 5,000 images for testing. For each image  $\mathbf{x}_r$ , we perform hair segmentation using off-the-shelf model, and calculate  $\mathbf{s}_r$  for the hair region.

**Fake Domain.** There are multiple ways to model and render virtual hair models. From coarse to fine, typical virtual hair models range from a single rigid shape, coarse polygon strips representing detached hair wisps, to large amount of thin hair fibers that mimic real-world hair behaviors. Due to various granularity of the geometry, the structural representation is hardly shared with each other or real hair images. In our experiments, all the hair models we used are polygon strips based considering this type of hair model is widely adopted in real-time scenarios for it is efficient to render and flexible to be animated. To generate  $\mathbf{s}_f$  for a given hair model  $\mathbf{h}$  and specified camera parameters  $\mathbf{c}$ , we use smoothly varying color gradient as texture to render  $\mathbf{h}$  into a color image that embeds the structure information of the hair geometry, such that  $\mathbf{s}_f = S_f(\mathbf{h}, \mathbf{c})$ . As for  $\mathbf{x}_f$ , we use traditional graphic rendering pipeline to render  $\mathbf{h}$  with a

uniform appearance color and simple diffuse shading, so that the final synthetic renderings have consistent appearance that can be easily disentangled without any extra condition, and keep all necessary structural information to verify the effectiveness of the encoding step. Example pairs are shown in Fig.2a.

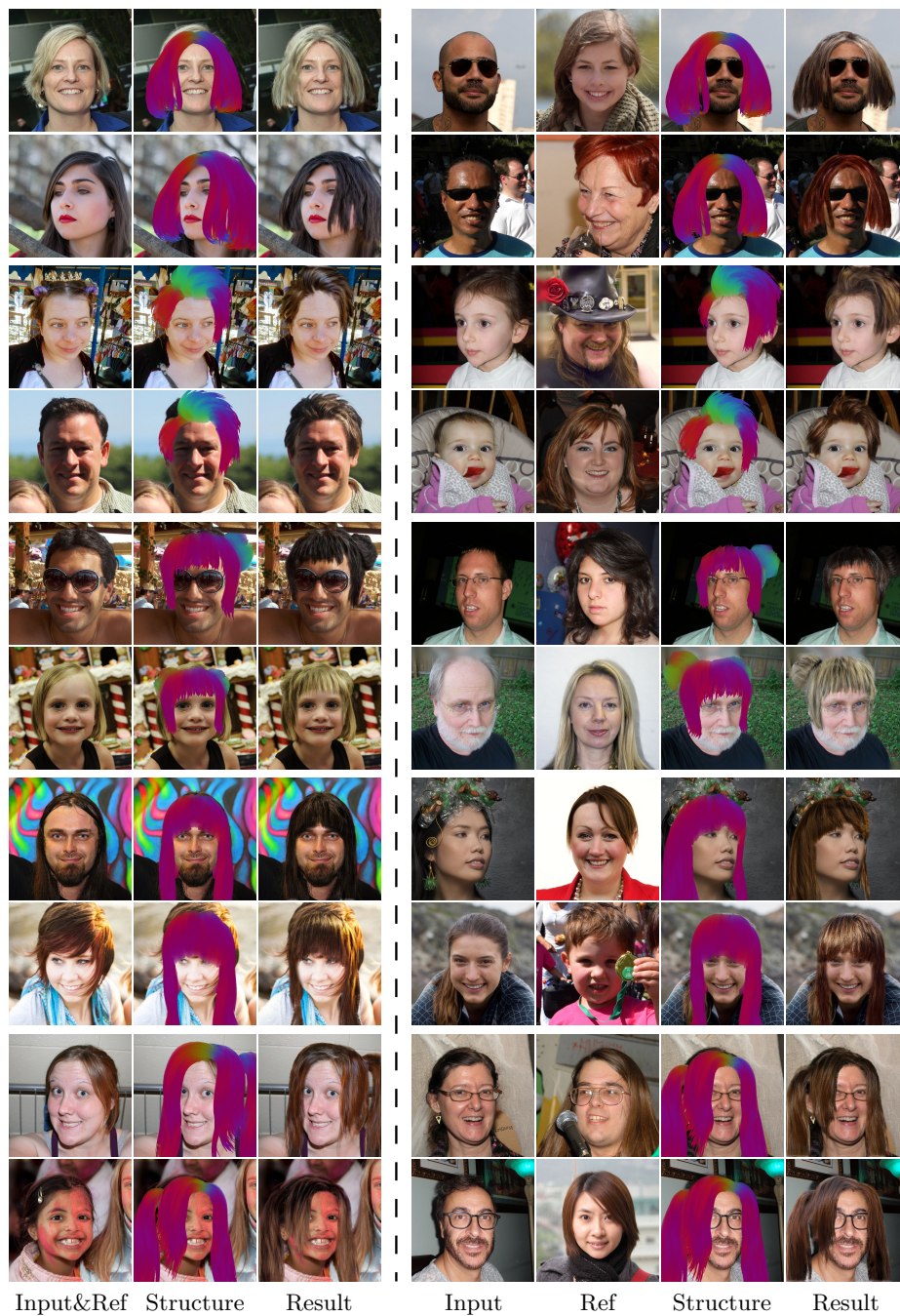
For the 3D hair used for fake data pairs, we create five models (leftmost column in Fig.2), including the middle hairstyle  $\mathbf{h}_m$ , the short punky hairstyle  $\mathbf{h}_s$ , the short hairstyle with two buns  $\mathbf{h}_b$ , the long hairstyle  $\mathbf{h}_l$ , and the long hairstyle with twin tails  $\mathbf{h}_t$ . The first four models are used for training, and the twin-tail hair model  $\mathbf{h}_t$  is used to evaluate the generalization capability of the network, for the network has never seen it. All these models consist of 10 to 50 polygon strips, which is sparse enough for real-time applications. We use the same training set from real domain to form training pairs. Each image is overlaid by one of the four 3D hair models according to the head position and pose. Then the image with fake hair model is used to generate  $\mathbf{x}_f$  through rendering the hair model with simple diffuse shading, and  $\mathbf{s}_f$  by exporting color textures that encodes surface tangent of the mesh. We strictly use the same shading parameters, including lighting and color, to enforce a uniform appearance of hair that can be easily disentangled by the networks.

## 4.2 Implementation Details

We apply a two-stage learning strategy. During the first stage, all networks are trained jointly following Eq.7 for the single-image renderer  $R_n$ . After that, we temporally fine-tune the decoder  $G_r$  of the real branch, to achieve temporally-smooth renderer  $R_n^t$ , by introducing the additional temporal condition as detailed in Sec.3.4. To make the networks of both stages consistent, we keep the same condition input dimensions, including appearance and temporal, but set the temporal condition to zero during the first stage. During the second stage we set it to zero with 50% of chance. The network architecture discussed in Sec.3 is implemented using PyTorch. We adopt Adam solver with a learning rate set to 0.0001 for the first stage, and 0.00001 for the fine-tuning stage. The training resolution of all images is  $512 \times 512$ , with the mini-batch size set to 4. For the loss functions, weights  $\lambda_p$ ,  $\lambda_s$ , and  $\lambda_g$  are set to 10, 1, and 1, respectively. All experiments are conducted on a workstation with 4 Nvidia Tesla P100 GPUs.

## 4.3 Qualitative Results

We present visual hair rendering results from two settings in Fig.3. The left three columns in Fig.3 show that the reference image  $\mathbf{r}$  is the same as  $\mathbf{x}_r$ . By applying a hair model, we can modify human hair shape but keep the original hair appearance and orientation. The right four columns show that the reference image is different from  $\mathbf{x}_r$ , therefore, both structure and appearance of hair from  $\mathbf{x}_r$  can be changed at the same time to render the hair with new style. These flexible applications demonstrate that our method can be easily applied to modify hair and generate novel high-quality hair images.



**Fig. 3. Selected results** for the hair models used in this study (2 examples per model). We visualize examples where the input and the reference image are the same (left), and the input and the reference are different images (right). In the former case the method copies appearance from another image.

#### 4.4 Comparison Results

To the best of our knowledge, there is no previous work that tackles the problem of neural hair rendering; thus, a direct comparison is not feasible. However, in light of our methods aim to bridge two different domains without ground-truth image pairs, which is related to unsupervised image translation, we compare our network with state-of-the-art unpaired image translation studies. It is important to stress that although our hair rendering translation falls into the range of image translation problems, there exist fundamental differences compared to the generic unpaired image translation formulations for the following two reasons.

First and foremost, compared with translation between two domains, such as painting styles, or seasons/times of the day, which have roughly the same amount of images for two domains and enough representative training images can be sampled to provide nearly-uniform domain coverage, our real/fake domains have dramatically different sizes—it is easy to collect a huge amount of real human portrait photos with diverse hairstyles to form the real domain. Unfortunately, for the fake domain, it is impossible to reach the same variety since it would require manually designing *every* possible hair shape and appearance to describe the distribution of the whole domain of rendered fake hair. Therefore, we focus on a realistic assumption that only a limited set of such models are available for training and testing, such that we use four 3D models for training and one for testing, which is far from being able to produce variety in the fake domain.

Second, as a deterministic process, hair rendering should be conditioned strictly on both geometric shape and chromatic appearance, which can be hardly achieved with unconditioned image translation frameworks.

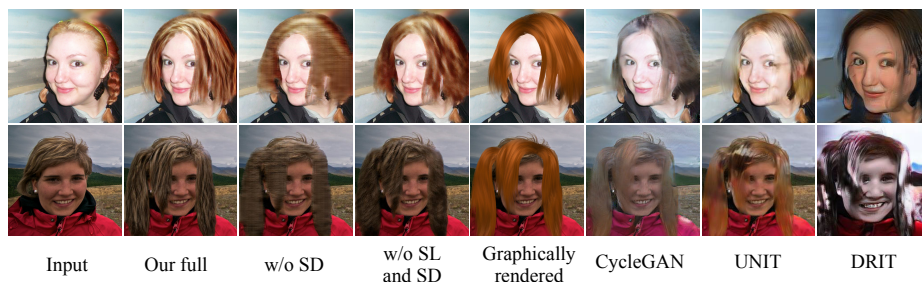
With those differences bearing in mind, we show the comparison between our method and three unpaired image translation studies, including CycleGAN [58], DRIT [26], and UNIT [30]. For the training of these methods, we use the same sets of images,  $\mathbf{x}_r$  and  $\mathbf{x}_f$ , for both real and fake domains, and the default hyperparameters reported by the original papers. Additionally, we compare with the images generated by traditional graphic rendering pipeline. We denote the method as **Graphic Renderer**. Finally, we report two ablation studies to evaluate the soundness of the network and the importance of each step: 1) we first remove the structural discriminator (termed as **w/o SD**); 2) we then additionally remove the shared latent space (termed as **w/o SL and SD**).

**Quantitative comparison.** For quantitative evaluation, we adopt FID (Frchet Inception Distance) [14] to measure the distribution distance between two domains. Moreover, inspired by the evaluation protocol from existing work [2,52], we apply a pre-trained hair segmentation model [46] on the generated images to get hair mask, and compare it with the ground truth hair mask. Intuitively, the segmentation model should predict the hair mask that similar to the ground-truth for the realistic synthesized images. To measure the segmentation accuracy, we use both Intersection-over-Union (IoU) and pixel accuracy (Accuracy).

The quantitative results are reported in Tab.1. Our method significantly outperforms the state-of-the-art unpaired image translation works and graphic

**Table 1. Quantitative comparison results.** We compare our method against commonly adopted image-to-image translation frameworks, reporting Frchet Inception Distance (FID, lower the better), Intersection over Union (IoU, higher the better) and pixel accuracy (Accuracy, higher the better). Additionally we report ablation studies by first removing the structural discriminator (w/o SD) followed by removing both the structural discriminator and the shared latent space (w/o SL and SD).

Method	FID ↓	IoU(%) ↑	Accuracy(%) ↑
Graphic Renderer	98.62	55.77	86.17
CycleGAN [58]	107.11	46.46	84.06
UNIT [30]	116.79	30.89	84.27
DRIT [26]	174.39	30.69	65.80
w/o SL and SD	94.25	80.10	93.89
w/o SD	77.09	86.60	96.35
<b>Ours</b>	<b>57.08</b>	<b>86.74</b>	<b>96.45</b>



**Fig. 4. Visual comparisons.** We show selected visual comparisons against commonly adopted image-to-image translation methods as well as visualize ablation results. Our method synthesizes more realistic hair images compared to other approaches.

rendering approach by a large margin for all the three evaluation metrics. The low FID score proves our method can generate high-fidelity hair images that contain similar hair appearance distribution as images from the real domain. The high IoU and Accuracy demonstrate the ability of the network to minimize structure gap between real and fake domain so that the synthesized images can follow the manually designed structure. Furthermore, the ablation analysis in Tab.1 shows both shared encoder layers and the structural discriminator are essential parts of the network, for the shared encoder layers help the network to find a common latent space that embeds hair structural knowledge, while the structural discriminator forces the hair structure features to be domain invariant.

**Qualitative comparison.** The qualitative comparison of different methods is shown in Fig.4. It can be easily seen that our generated images have much higher quality than the synthesized images created by other state-of-the-art unpaired image translation methods, for they have clearer hair mask, follow hair appearance from reference images, maintain the structure from hair models, and look



**Fig. 5. Video results and comparisons.** Top row: the first image is the appearance reference image and others are continuous input frames; middle row: generated hair images with temporal conditioning; bottom row: generated hair images without temporal conditioning. We show two zoom-in hair regions for each result. By applying temporal conditioning, our model synthesizes hair images with consistent appearance, while not using temporal conditioning leads to hair appearance flickering as indicated by blue and green boxes. *Click the image to play the video results and comparisons.*

like natural hair. Compared with the ablation methods (Fig.4c and d), our full method (Fig.4b) can follow the appearance from reference images (Fig.4a) by generating hair with similar orientation.

We also show the importance of temporal conditioning (Sec.3.4) in Fig.5. The temporal conditioning helps us generate consistent and smooth video results, for hair appearance and orientation are similar between continuous frames. Without temporal conditioning, the hair texture could be different between frames, as indicated by blue and green boxes, which may result in flickering for the synthesized video. Please refer to the supplementary video for more examples.

## 5 Conclusions

We propose a neural-based rendering pipeline for general virtual 3D hair models. The key idea of our method is that instead of enforcing model-level representation consistency to enable supervised paired training, we relax the strict requirements on the model and adopt a unsupervised image translation framework. To bridge the gap between real and fake domains, we construct a shared latent space to encode a common structure feature space for both domains, even if their inputs are dramatically different. In this way, we can encode a virtual hair model into such a structure feature, and switch it into the real generator to produce realistic rendering. The conditional real generator not only allow flexible condition of hair appearance, but can also be used to introduce an extra temporal conditioning to generate smooth sequential results.

## References

1. Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., Krishnan, D.: Unsupervised pixel-level domain adaptation with generative adversarial networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3722–3731 (2017)
2. Chen, Q., Koltun, V.: Photographic image synthesis with cascaded refinement networks. In: Proceedings of the IEEE international conference on computer vision. pp. 1511–1520 (2017)
3. Chen, T.Q., Schmidt, M.: Fast patch-based style transfer of arbitrary style. arXiv preprint arXiv:1612.04337 (2016)
4. Chen, Y.H., Chen, W.Y., Chen, Y.T., Tsai, B.C., Frank Wang, Y.C., Sun, M.: No more discrimination: Cross city adaptation of road scene segmenters. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1992–2001 (2017)
5. d’Eon, E., François, G., Hill, M., Letteri, J., Aubry, J.: An energy-conserving hair reflectance model. *Comput. Graph. Forum* **30**(4), 1181–1187 (2011)
6. Dundar, A., Liu, M.Y., Wang, T.C., Zedlewski, J., Kautz, J.: Domain stylization: A strong, simple baseline for synthetic to real image domain adaptation. arXiv preprint arXiv:1807.09384 (2018)
7. Fernando, B., Habrard, A., Sebban, M., Tuytelaars, T.: Unsupervised visual domain adaptation using subspace alignment. In: Proceedings of the IEEE international conference on computer vision. pp. 2960–2967 (2013)
8. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. arXiv preprint arXiv:1409.7495 (2014)
9. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.: Domain-adversarial training of neural networks. *The Journal of Machine Learning Research* **17**(1), 2096–2030 (2016)
10. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2414–2423 (2016)
11. Gong, B., Shi, Y., Sha, F., Grauman, K.: Geodesic flow kernel for unsupervised domain adaptation. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 2066–2073. IEEE (2012)
12. Gopalan, R., Li, R., Chellappa, R.: Domain adaptation for object recognition: An unsupervised approach. In: 2011 international conference on computer vision. pp. 999–1006. IEEE (2011)
13. Hertzmann, A., Jacobs, C.E., Oliver, N., Curless, B., Salesin, D.H.: Image analogies. In: Proceedings of the 28th annual conference on Computer graphics and interactive techniques. pp. 327–340 (2001)
14. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: Advances in Neural Information Processing Systems. pp. 6626–6637 (2017)
15. Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A.A., Darrell, T.: Cycada: Cycle-consistent adversarial domain adaptation. arXiv preprint arXiv:1711.03213 (2017)
16. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1501–1510 (2017)

17. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 172–189 (2018)
18. Isola, P., Zhu, J., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR. pp. 5967–5976 (2017)
19. Jo, Y., Park, J.: SC-FEGAN: Face editing generative adversarial network with user’s sketch and color. In: IEEE International Conference on Computer Vision, ICCV 2019. pp. 1745–1753 (2019)
20. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: European conference on computer vision. pp. 694–711. Springer (2016)
21. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR. pp. 4401–4410 (2019)
22. Kim, T., Cha, M., Kim, H., Lee, J.K., Kim, J.: Learning to discover cross-domain relations with generative adversarial networks. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 1857–1865. JMLR. org (2017)
23. Kulis, B., Saenko, K., Darrell, T.: What you saw is not what you get: Domain adaptation using asymmetric kernel transforms. In: CVPR 2011. pp. 1785–1792. IEEE (2011)
24. Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., et al.: Photo-realistic single image super-resolution using a generative adversarial network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4681–4690 (2017)
25. Lee, C.H., Liu, Z., Wu, L., Luo, P.: Maskgan: Towards diverse and interactive facial image manipulation (2019)
26. Lee, H., Tseng, H., Huang, J., Singh, M., Yang, M.: Diverse image-to-image translation via disentangled representations. In: European Conference on Computer Vision, ECCV. vol. 11205, pp. 36–52 (2018)
27. Li, C., Wand, M.: Precomputed real-time texture synthesis with markovian generative adversarial networks. In: European conference on computer vision. pp. 702–716. Springer (2016)
28. Li, Y., Wang, N., Liu, J., Hou, X.: Demystifying neural style transfer. arXiv preprint arXiv:1701.01036 (2017)
29. Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.H.: Diversified texture synthesis with feed-forward networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3920–3928 (2017)
30. Liu, M., Breuel, T., Kautz, J.: Unsupervised image-to-image translation networks. In: Annual Conference on Neural Information Processing Systems, NeurIPS. pp. 700–708 (2017)
31. Liu, M.Y., Huang, X., Mallya, A., Karras, T., Aila, T., Lehtinen, J., Kautz, J.: Few-shot unsupervised image-to-image translation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 10551–10560 (2019)
32. Liu, M.Y., Tuzel, O.: Coupled generative adversarial networks. In: Advances in neural information processing systems. pp. 469–477 (2016)
33. Marschner, S.R., Jensen, H.W., Cammarano, M., Worley, S., Hanrahan, P.: Light scattering from human hair fibers. *ACM Trans. Graph.* **22**(3), 780–791 (2003)
34. Moon, J.T., Marschner, S.R.: Simulating multiple scattering in hair using a photon mapping approach. *ACM Trans. Graph.* **25**(3), 1067–1074 (2006)



35. Moon, J.T., Walter, B., Marschner, S.: Efficient multiple scattering in hair using spherical harmonics. *ACM Trans. Graph.* **27**(3), 31 (2008)
36. Paris, S., Chang, W., Kozhushnyan, O.I., Jarosz, W., Matusik, W., Zwicker, M., Durand, F.: Hair photobooth: geometric and photometric acquisition of real hairstyles. *ACM Trans. Graph.* **27**(3), 30 (2008)
37. Park, T., Liu, M., Wang, T., Zhu, J.: Semantic image synthesis with spatially-adaptive normalization. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. pp. 2337–2346 (2019)
38. Qiu, H., Wang, C., Zhu, H., Zhu, X., Gu, J., Han, X.: Two-phase hair image synthesis by self-enhancing generative model. *Comput. Graph. Forum* **38**(7), 403–412 (2019)
39. Ren, Z., Zhou, K., Li, T., Hua, W., Guo, B.: Interactive hair rendering under environment lighting. *ACM Trans. Graph.* **29**(4), 55:1–55:8 (2010)
40. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.S., Berg, A.C., Li, F.: Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* **115**(3), 211–252 (2015)
41. Sadeghi, I., Pritchett, H., Jensen, H.W., Tamstorf, R.: An artist friendly hair shading system. *ACM Trans. Graph.* **29**(4), 56:1–56:10 (2010)
42. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting visual category models to new domains. In: *European conference on computer vision*. pp. 213–226. Springer (2010)
43. Sangkloy, P., Lu, J., Fang, C., Yu, F., Hays, J.: Scribbler: Controlling deep image synthesis with sketch and color. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5400–5409 (2017)
44. Shrivastava, A., Pfister, T., Tuzel, O., Susskind, J., Wang, W., Webb, R.: Learning from simulated and unsupervised images through adversarial training. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2107–2116 (2017)
45. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: *International Conference on Learning Representations, ICLR* (2015)
46. Svanera, M., Muhammad, U.R., Leonardi, R., Benini, S.: Figaro, hair detection and segmentation in the wild. In: *2016 IEEE International Conference on Image Processing (ICIP)*. pp. 933–937. IEEE (2016)
47. Taigman, Y., Polyak, A., Wolf, L.: Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200* (2016)
48. Tsai, Y.H., Hung, W.C., Schulter, S., Sohn, K., Yang, M.H., Chandraker, M.: Learning to adapt structured output space for semantic segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 7472–7481 (2018)
49. Tzeng, E., Hoffman, J., Saenko, K., Darrell, T.: Adversarial discriminative domain adaptation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 7167–7176 (2017)
50. Tzeng, E., Hoffman, J., Zhang, N., Saenko, K., Darrell, T.: Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474* (2014)
51. Ulyanov, D., Lebedev, V., Vedaldi, A., Lempitsky, V.S.: Texture networks: Feed-forward synthesis of textures and stylized images. In: *ICML*. vol. 1, p. 4 (2016)
52. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. In:

- Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8798–8807 (2018)
53. Ward, K., Bertails, F., Kim, T., Marschner, S.R., Cini, M., Lin, M.C.: A survey on hair modeling: Styling, simulation, and rendering. *IEEE Trans. Vis. Comput. Graph.* **13**(2), 213–234 (2007)
  54. Wei, L., Hu, L., Kim, V.G., Yumer, E., Li, H.: Real-time hair rendering using sequential adversarial networks. In: *European Conference on Computer Vision, ECCV*. vol. 11208, pp. 105–122 (2018)
  55. Xu, K., Ma, L., Ren, B., Wang, R., Hu, S.: Interactive hair rendering and appearance editing under environment lighting. *ACM Trans. Graph.* **30**(6), 173 (2011)
  56. Yan, L., Tseng, C., Jensen, H.W., Ramamoorthi, R.: Physically-accurate fur reflectance: modeling, measurement and rendering. *ACM Trans. Graph.* **34**(6), 185:1–185:13 (2015)
  57. Yi, Z., Zhang, H., Tan, P., Gong, M.: Dualgan: Unsupervised dual learning for image-to-image translation. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2849–2857 (2017)
  58. Zhu, J., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: *IEEE International Conference on Computer Vision, ICCV*. pp. 2242–2251 (2017)
  59. Zinke, A., Yuksel, C., Weber, A., Keyser, J.: Dual scattering approximation for fast multiple scattering in hair. *ACM Trans. Graph.* **27**(3), 32 (2008)