

3DAvatarGAN: Bridging Domains for Personalized Editable Avatars

Rameen Abdal^{†1} Hsin-Ying Lee² Peihao Zhu^{†1} Menglei Chai² Aliaksandr Siarohin²

Peter Wonka¹ Sergey Tulyakov²

¹KAUST ²Snap Inc.

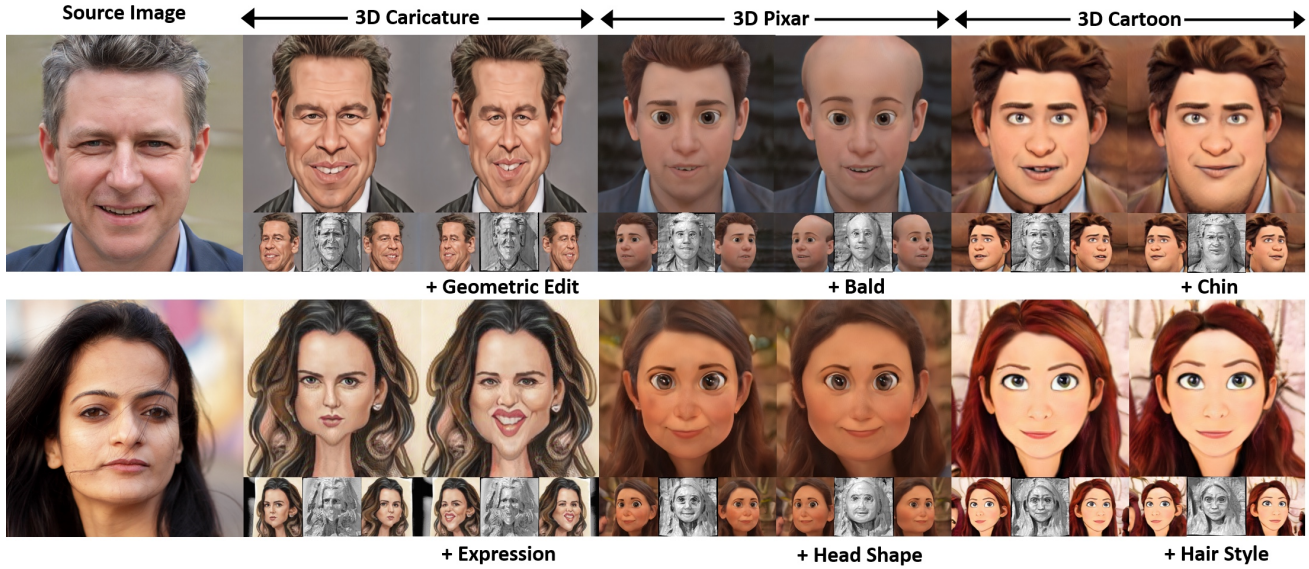


Figure 1. **Editable 3D avatars.** We present 3DAvatarGAN, a 3D GAN able to produce and edit personalized 3D avatars from a single photograph (real or generated). Our method distills information from a 2D-GAN trained on 2D artistic datasets like Caricatures, Pixar toons, Cartoons, Comics *etc.* and requires no camera annotations.

Abstract

Modern 3D-GANs synthesize geometry and texture by training on large-scale datasets with a consistent structure. Training such models on stylized, artistic data, with often unknown, highly variable geometry, and camera information has not yet been shown possible. Can we train a 3D GAN on such artistic data, while maintaining multi-view consistency and texture quality? To this end, we propose an adaptation framework, where the source domain is a pre-trained 3D-GAN, while the target domain is a 2D-GAN trained on artistic datasets. We, then, distill the knowledge from a 2D generator to the source 3D generator. To do that, we first propose an optimization-based method to align the distributions of camera parameters across domains. Second, we propose regularizations necessary to learn high-quality texture, while avoiding degenerate geometric solutions, such as flat shapes. Third, we show

a deformation-based technique for modeling exaggerated geometry of artistic domains, enabling—as a byproduct—personalized geometric editing. Finally, we propose a novel inversion method for 3D-GANs linking the latent spaces of the source and the target domains. Our contributions—for the first time—allow for the generation, editing, and animation of personalized artistic 3D avatars on artistic datasets. Project Page: <https://rameenabdal.github.io/3DAvatarGAN>

1. Introduction

Photo-realistic portrait face generation is an iconic application demonstrating the capability of generative models especially GANs [29, 31, 32]. A recent development has witnessed an advancement from straightforwardly synthesizing 2D images to learning 3D structures without 3D supervision, referred to as 3D-GANs [10, 43, 58, 67]. Such training

[†] Part of the work was done during an internship at Snap Inc.

is feasible with the datasets containing objects with highly consistent geometry, enabling a 3D-GAN to learn a distribution of shapes and textures. In contrast, artistically stylized datasets [26, 68] have arbitrary exaggerations of both geometry and texture, for example, the nose, cheeks, and eyes can be arbitrarily drawn, depending on the style of the artist as well as on the features of the subject, see Fig. 1. Training a 3D-GAN on such data becomes problematic due to the challenge of learning such an arbitrary distribution of geometry and texture. In our experiments (Sec. 5.1), 3D-GANs [10] generate flat geometry and become 2D-GANs essentially. A natural question arises, whether a 3D-GAN can synthesize consistent novel views of images belonging to artistically stylized domains, such as the ones in Fig. 1.

In this work, we propose a domain-adaption framework that allows us to answer the question positively. Specifically, we fine-tune a pre-trained 3D-GAN using a 2D-GAN trained on a target domain. Despite being well explored for 2D-GANs [26, 68], existing domain adaptation techniques are not directly applicable to 3D-GANs, due to the nature of 3D data and characteristics of 3D generators.

The geometry and texture of stylized 2D datasets can be arbitrarily exaggerated depending on the context, artist, and production requirements. Due to this, no reliable way to estimate camera parameters for each image exists, whether using an off-the-shelf pose detector [75] or a manual labeling effort. To enable the training of 3D-GANs on such challenging datasets, we propose three contributions. ① An optimization-based method to align distributions of camera parameters between domains. ② Texture, depth, and geometry regularizations to avoid degenerate, flat solutions and ensure high visual quality. Furthermore, we redesign the discriminator training to make it compatible with our task. We then propose ③ a *Thin Plate Spline (TPS)* 3D deformation module operating on a tri-plane representation to allow for certain large and sometimes extreme geometric deformations, which are so typical in artistic domains.

The proposed adaptation framework enables the training of 3D-GANs on complex and challenging artistic data. The previous success of domain adaptation in 2D-GANs unleashed a number of exciting applications in the content creation area [26, 68]. Given a single image such methods first find a latent code corresponding to it using GAN inversion, followed by latent editing producing the desired effect in the image space. Compared to 2D-GANs, the latent space of 3D-GANs is more entangled, making it more challenging to link the latent spaces between domains, rendering the existing inversion and editing techniques not directly applicable. Hence, we take a step further and explore the use of our approach to 3D artistic avatar generation and editing. Our final contribution to enable such applications is ④ a new inversion method for coupled 3D-GANs.

In summary, the proposed domain-adaption framework

allows us to train 3D-GANs on challenging artistic datasets with exaggerated geometry and texture. We call our method 3DAvatarGAN as it—for the first time—offers generation, editing, and animation of personalized stylized, artistic avatars obtained from a single image. Our results (See Sec. 5.2) show the high-quality 3D avatars possible by our method compared to the naive fine-tuning.

2. Related Work

GANs and Semantic Image Editing. Generative adversarial Networks (GANs) [20, 50] are one popular type of generative model, especially for smaller high-quality datasets such as FFHQ [33], AFHQ [14], and LSUN objects [70]. For these datasets, StyleGAN [29, 31, 33] can be considered as the current state-of-the-art GAN [28, 29, 31, 33, 34]. The disentangled latent space learned by StyleGAN has been shown to exhibit semantic properties conducive to semantic image editing [1, 3, 16, 23, 37, 46, 54, 59, 65]. CLIP [49] based image editing [2, 17, 46] and domain transfer [15, 73] are another set of works enabled by StyleGAN.

GAN Inversion. Algorithms to project existing images into a GAN latent space are a prerequisite for GAN-based image editing. There are mainly two types of methods to enable such a projection: optimization-based methods [1, 13, 60, 74] and encoder-based methods [5, 7, 51, 61, 72]. On top of both streams of methods, the generator weights can be further modified after obtaining initial inversion results [52].

Learning 3D-GANs with 2D Data. Previously, some approaches attempt to extract 3D structure from pre-trained 2D-GANs [44, 55]. Recently, inspired by Neural Radiance Field (NeRF) [9, 38, 45, 71], novel GAN architectures have been proposed to combine implicit or explicit 3D representations with neural rendering techniques [11, 12, 21, 40, 41, 43, 53, 56, 58, 66, 67]. In our work, we build on EG3D [11] which has current state-of-the-art results for human faces trained on the FFHQ dataset.

Avatars and GANs. To generate new results in an artistic domain (e.g. anime or cartoons), a promising technique is to fine-tune an existing GAN pre-trained on photographs, e.g. [47, 57, 63]. Data augmentation and freezing lower layers of the discriminator are useful tools when fine-tuning a 2D-GAN [29, 39]. One branch of methods [18, 46, 73] investigates domain adaptation if only a few examples or only text descriptions are available. While others focus on matching the distribution of artistic datasets with diverse shapes and styles. Our work also falls in this domain. Among previous efforts, StyleCariGAN [26] proposes invertible modules in the generator to train and generate caricatures from real images. DualStyleGAN [68] learns two mapping networks in StyleGAN to control the style and structure of the new domain. Some works are trained on 3D data or require heavy labeling/engineering [22, 27, 69] and use 3D morphable models to map 2D images of carica-

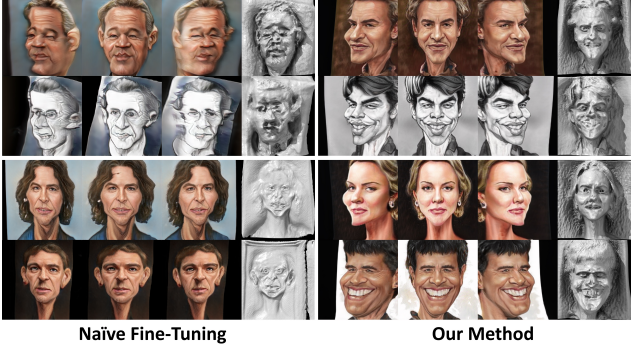


Figure 2. **Comparison with naive fine-tuning.** Comparison of generated 3D avatars with a naïvely fine-tuned generator G_{base} (left sub-figures) versus our generator G_t (right sub-figures). The corresponding sub-figures show comparisons in terms of texture quality (top two rows) and geometry (bottom two rows). See Sec. 5.1 for details.

tures to 3D models. However, such models fail to model the hair, teeth, neck, and clothes and suffer in texture quality. In this work, we are the first to tackle the problem of domain adaption of 3D-GANs and to produce fully controllable 3D Avatars. We employ 2D to 3D domain adaptation and distillation and make use of synthetic 2D data from StyleCariGAN [26] and DualStyleGAN [68].

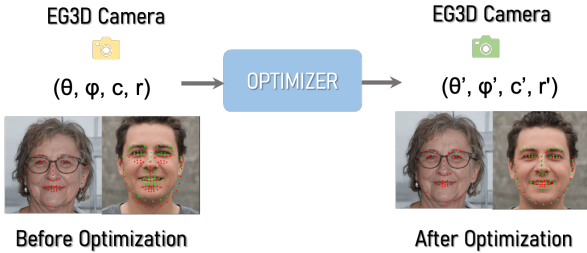


Figure 3. **Illustration of camera alignment.** We define an optimization algorithm to determine the camera parameters of the new domain *i.e.* avatars.

3. Domain Adaptation for 3D-GANs

The goal of domain adaptation for 3D-GANs is to adapt (both texture and geometry) to a particular style defined by a 2D dataset (Caricature, Anime, Pixar toons, Comic, and Cartoons [25, 26, 68] in our case). In contrast to 2D-StyleGAN-based fine-tuning methods that are conceptually simpler [30, 47], fine-tuning a 3D-GAN on 2D data introduces challenges in addition to domain differences, especially on maintaining the texture quality while preserving the geometry. Moreover, for these datasets, there is no explicit shape and camera information. We define the domain adaptation task as follows: Given a prior 3D-GAN *i.e.* EG3D (G_s) of source domain (T_s), we aim to produce

a 3D Avatar GAN (G_t) of the target domain (T_t) while maintaining the semantic, style, and geometric properties of G_s , and at the same time preserving the identity of the subject between the domains ($T_s \leftrightarrow T_t$). Refer to Fig. 5 for the pipeline figure. We represent G_{2D} as a teacher 2D-GAN used for knowledge distillation fine-tuned on the above datasets. Note that as T_t is not assumed to contain camera parameter annotations, the training scheme must suppress artifacts such as low-quality texture under different views and flat geometry (See Fig. 2). In the following, we discuss the details of our method.

3.1. How to align the cameras?

Selecting appropriate ranges for camera parameters is of paramount importance for high-fidelity geometry and texture detail. Typically, such parameters are empirically estimated, directly computed from the dataset using an off-the-shelf pose detector [10], or learned during training [8]. In domains we aim to bridge, such as caricatures for which a 3D model may not even exist, directly estimating the camera distribution is problematic and, hence, is not assumed by our method. Instead, we find it essential to ensure that the camera parameter distribution is consistent across the source and target domains. For the target domain, we use StyleGAN2 trained on FFHQ, fine-tuned on artistic datasets [26, 68]. Assuming that the intrinsic parameters of all the cameras are the same, we aim to match the distribution of extrinsic camera parameters of G_s and G_{2D} and train our final G_t using it (see illustration in Fig. 3). To this end, we define an optimization-based method to match the sought distributions. The first step is to identify a canonical pose image in G_{2D} , where the yaw, pitch, and roll parameters are zero. According to Karras et al., [32], the image corresponding to the mean latent code satisfies this property. Let θ, ϕ be the camera Euler angles in a spherical coordinate system, r, c be the radius of the sphere and camera lookat point, and M be a function that converts these parameters into the camera-to-world matrix. Let $I_s(w, \theta, \phi, c, r) = G_s(w, M(\theta, \phi, c, r))$ and $I_{2D}(w) = G_{2D}(w)$ represent an arbitrary image generated by G_s and G_{2D} , respectively, given the w code variable. Let k_d be the face key-points detected by the detector K_d [75], then

$$(c', r') := \arg \min_{(c, r)} L_{kd}(I_s(w'_{\text{avg}}, 0, 0, c, r), I_{2D}(w_{\text{avg}})), \quad (1)$$

where $L_{kd}(I_1, I_2) = \|k_d(I_1) - k_d(I_2)\|_1$ and w_{avg} and w'_{avg} are the mean w latent codes of G_{2D} and G_s , respectively. In our results, r' is determined to be 2.7 and c' is approximately $[0.0, 0.05, 0.17]$. The next step is to determine a safe range of the θ and ϕ parameters. Following prior works, StyleFlow [3] and FreeStyleGAN [36] (see Fig. 5 of the paper), we set these parameters as $\theta' \in [-0.45, 0.45]$ and $\phi' \in [-0.35, 0.35]$ in radians.

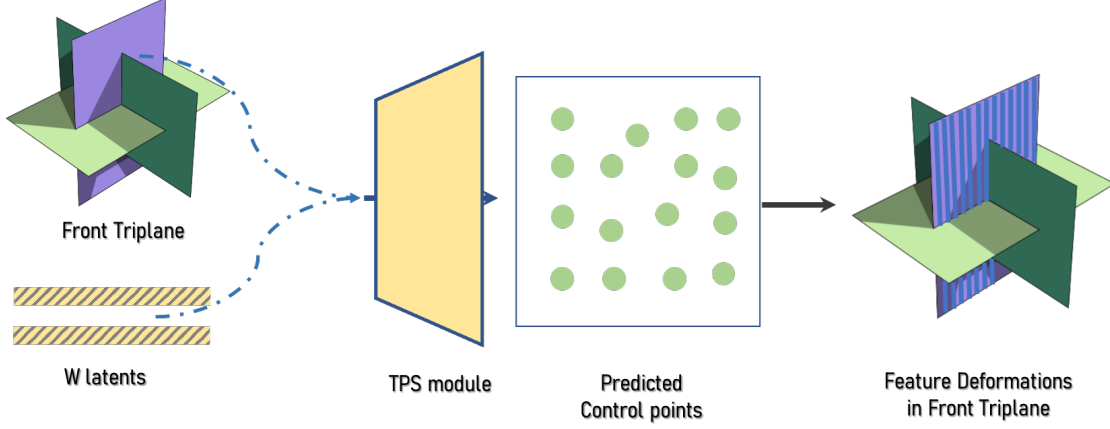


Figure 4. **Illustration of TPS module.** Our *TPS* module takes W latents and the front triplane as inputs and outputs deformed control points that deform the front triplane.

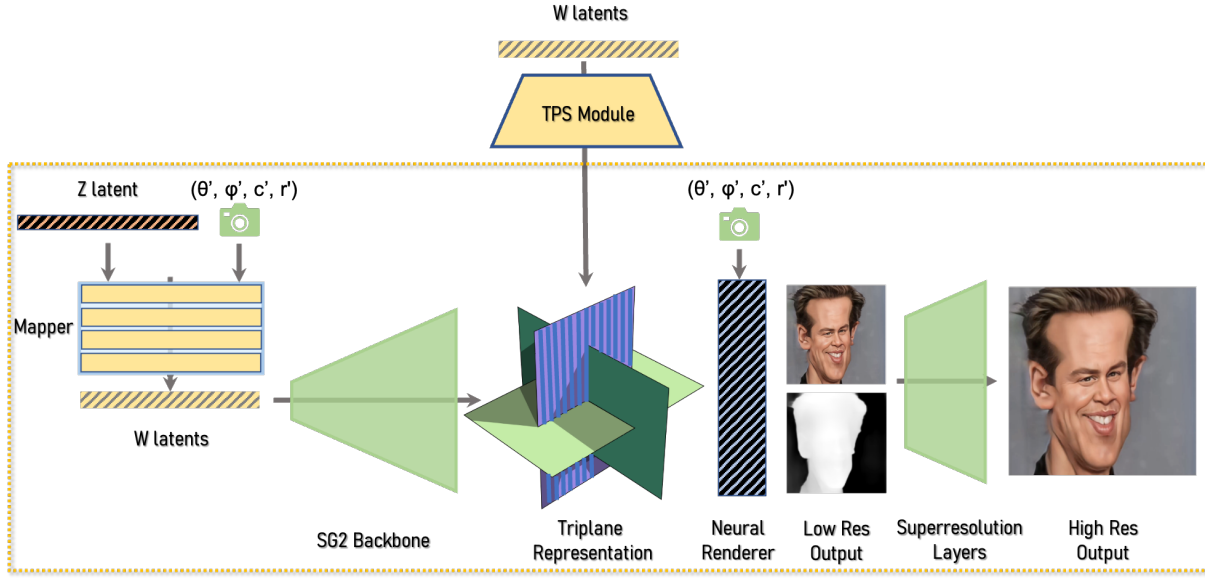


Figure 5. **Fine tuned EG3D (G_t) pipeline with *TPS* module.** We design a framework to fine-tune a 3D-GAN *i.e.* EG3D (dotted yellow box) and adapt to an artistic domain that does not stem from a consistent 3D model. In addition, our framework modifies EG3D architecture to include a novel *TPS* module for geometric editing.

3.2. What loss functions and regularizers to use?

Next, although the camera systems are aligned, the given dataset may not stem from a consistent 3D model, *e.g.*, in the case of caricatures or cartoons. This entices the generator G_t to converge to an easier degenerate solution with a flat geometry. Hence, to benefit from the geometric prior of G_s , another important step is to design the loss functions and regularizers for a selected set of parameters to update in G_t . Next, we discuss these design choices:

Loss Functions. To ensure texture quality and diversity, we resort to the adversarial loss used to fine-tune GANs as our main loss function. We use the standard non-saturating

loss to train the generator and discriminator networks used in EG3D [11]. We also perform lazy density regularization to ensure consistency of the density values in the final fine-tuned model G_t .

Texture Regularization. Since the texture can be entangled with the geometry information, determining which layers to update is important. To make use of the fine-style information encoded in later layers, it is essential to update the *tRGB* layer parameters (outputting tri-plane features) before the neural rendering stage. *tRGB* are convolutional layers that transform feature maps to 3 channels at each resolution (96 channels in triplanes). Moreover, since the net-



Figure 6. **Domain adaptation.** Domain adaptation results of images from source domain T_s (top row in each sub-figure) to target domain T_t . Rows two to five show corresponding 3D avatar results from different viewpoints.



Figure 7. **Real to avatar results.** Our framework takes a single image as input and generates editable 3D Avatars. Also see Fig. 8

work has to adapt to a color distribution of T_t , it is essential to update the decoder (*MLP* layers) of the neural rendering pipeline as well. Given the EG3D architecture, we also update the super-resolution layer parameters to ensure the coherency between the low-resolution and high-resolution outputs seen by the discriminator D .

Geometry Regularization. In order to allow the network to learn the structure distribution of T_t and at the same time ensure properties of \mathcal{W} and \mathcal{S} latent spaces are preserved, we update the earlier layers with regularization. This also encourages the latent spaces of T_s and T_t to be easily linked. Essentially, we update the deviation parameter Δs from the s activations of the \mathcal{S} space [65]. The s activations are predicted by $A(w)$, where A is the learned affine function in EG3D. The s activations scale the kernels of a particular layer. In order to preserve the identity as well as geometry such that the optimization of Δs does not deviate too far

away from the original domain T_s , we introduce a regularizer given by

$$R(\Delta s) := \|\Delta s\|_1. \quad (2)$$

Note that we apply $R(\Delta s)$ regularization in a lazy manner, *i.e.*, with density regularization. Interestingly, after training, we can interpolate between s and $s + \Delta s$ parameters to interpolate between the geometries of samples in T_s and T_t (See Fig. 9).

Depth Regularization. Next, we observe that even though the above design choice produces better geometry for T_t , some samples from G_t can still lead to flatter geometry, and it is hard to detect these cases. We found that the problem is related to the relative depth of the background to the foreground. To circumvent this problem, we use an additional regularization where we encourage the average background depth of G_t to be similar to G_s . Let S_b be a face background segmentation network [35]. We first compute the

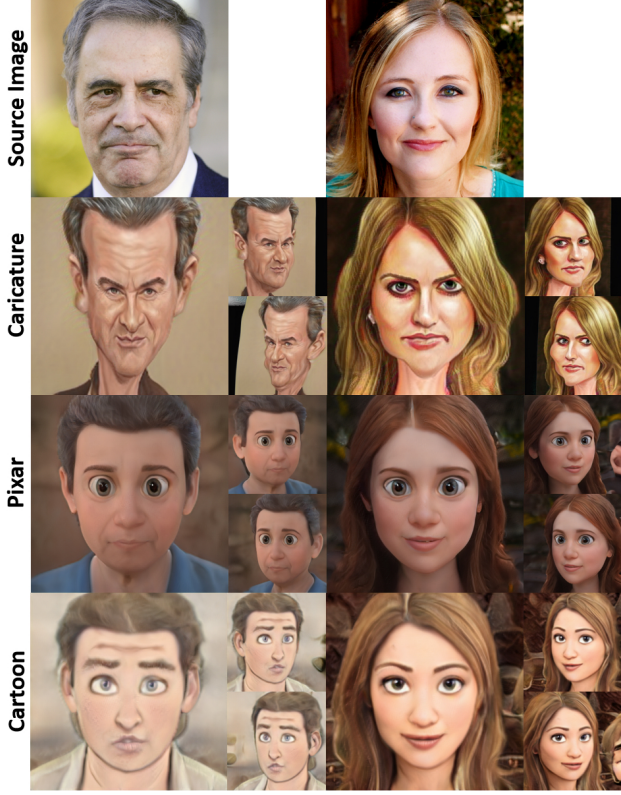


Figure 8. **3D avatars from real images.** Projection of real images on the 3D avatar generators.

average background depth of the samples given by G_s . This average depth is given by

$$a_d := \frac{1}{M} \sum_{n=1}^M \left(\frac{1}{N_n} \|D_n \odot S_b(I_n)\|_F^2 \right). \quad (3)$$

Here, D_n is the depth map of the image I_n sampled from G_s , \odot represents the *Hadamard* product, M is the number of the sampled images, and N_n is the number of background pixels in I_n . Finally, regularization is defined as:

$$R(D) := \|a_d \cdot J - (D_t \odot S_b(I_t))\|_F, \quad (4)$$

where D_t is the depth map of the image I_t sampled from G_t and J is the matrix of ones having the same spatial dimensions as D_t .

3.3. What discriminator to use?

Given that the data in T_s and T_t is not paired and T_t is not assumed to contain camera parameter annotations, the choice of the discriminator (D) used for this task is also a critical design choice. Essentially, we use the unconditional version of the dual discriminator proposed in EG3D, and hence, we do not condition the discriminator on the camera information. As a result, during the training, G_t generates

arbitrary images with pose using $M(\theta', \phi', c', r')$, and the discriminator discriminates these images using arbitrary images from T_t . We train the discriminator from scratch and in order to adapt $T_s \rightarrow T_t$, we use the StyleGAN-ADA [29] training scheme and use R1 regularization.

3.4. How to incorporate larger geometric deformations between domains?

While the regularizers are used to limit the geometric changes when adapting from T_s to T_t , modeling large geometric deformations, e.g., in the caricature dataset is another challenge. One choice to edit the geometry is to use the properties of tri-plane features learned by EG3D. We start out by analyzing these three planes in G_s . We observe that the frontal plane encodes most of the information required to render the final image. To quantify this, we sample images and depth maps from G_s and swap the front and the other planes from two random images. Then we compare the difference in *RGB* values of the images and the *Chamfer* distance of the depth maps. While swapping the frontal tri-planes, the final images are completely swapped, and the *Chamfer* distance changes by 80 ~ 90% matching the swapped image depth map. In the case of the other two planes, the *RGB* image is not much affected and the *Chamfer* distance of the depth maps is reduced by only 20 ~ 30% in most cases.

Given the analysis, we focus to manipulate the 2D front plane features to learn additional deformation or exaggerations. We learn a *TPS* (*Thin Plate Spline*) [64] network on top of the front plane. Our *TPS* network is conditioned both on the front plane features as well as the \mathcal{W} space to enable multiple transformations. The architecture of the module is similar to the standard StyleGAN2 layer with an *MLP* appended at the end to predict the control points that transform the features. Hence, as a byproduct, we also enable 3D-geometry editing guided by the learned latent space. We train this module separately after G_t has been trained. We find that joint training is unstable due to exploding gradients arising from the large domain gap between T_s and T_t in the initial stages. Formally, we define this transformation as:

$$T(w, f) := \Delta c, \quad (5)$$

where, w is the latent code, f is the front plane, and c are the control points.

Let c_1 be the initial control points producing an identity transformation, (c_1, c_2) be the control points corresponding to front planes (f_1, f_2) sampled using \mathcal{W} codes (w_1, w_2) , respectively, and (c'_1, c'_2) be points with (w_1, w_2) swapped in the *TPS* module. To regularize and encourage the module

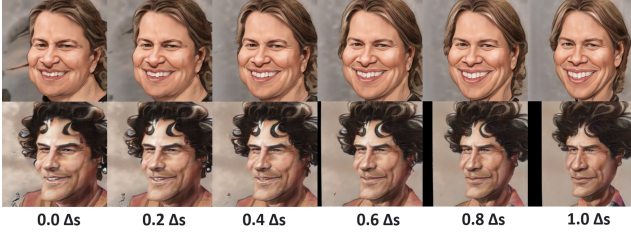


Figure 9. **Interpolation of Δs .** Geometric deformation using the interpolation of learned Δs parameters.

to learn different deformations, we have

$$R(T_1) := \alpha \sum_{n=1}^2 \|c_I - c_n\|_1 - \beta \|c_1 - c_2\|_1 - \sigma \|c'_1 - c'_2\|_1. \quad (6)$$

We use initial control point regularization to regularize large deviations in the control points which would otherwise explode. Additionally, to learn extreme exaggerations in T_t and ‘in expectation’, conform to the target distribution in the dataset, we add an additional loss term. Let $S(I)$ be the softmax output of the face segmentation network [35] given an image I and assuming that S generalizes to caricatures, then

$$R(T_2) := \|S(G_t(w)), S(I_t)\|_1 \quad (7)$$

Eq. 6, Eq. 7, and adversarial training loss are used to train the *TPS* module. We adopt gradient clipping to make sure that the training does not diverge. See the illustrations in Fig. 4 and Fig. 5.

4. Personalized Avatar Generation and Editing

Although 3D domain adaptation adapts $T_s \leftrightarrow T_t$, it is still a challenge to effectively link the latent spaces of G_s and G_t to generate personalized 3D avatars using a single photograph as the reference image. Particularly, the challenge arises due to the discrepancy in the coupled latent spaces when dealing with the projection of real photographs on 3D generators. Moreover, one would like to edit and animate these 3D avatars.

Projection. The task is to project a real image into the latent space of G_s , transfer the latent to G_t , and further optimize it to construct a 3D avatar. First, we use an optimization-based method to find the w code that minimizes the similarity between the generated and the real image in G_s . To achieve this, the first step is to align the cameras. We follow the steps mentioned in Sec. 3.1 for this step. Next, we use pixel-wise *MSE* loss and *LPIPS* loss to project the image into G_s [1]. Additionally, to preserve the identity of the subject, we use attribute classifiers *e.g.* caricature dataset [25] provides the coupled attribute information of real images and caricatures. We use such attribute classifier [25, 26] in a post-hoc manner as we notice that such networks can affect

the texture in the target domain and could degenerate to narrow style outputs if applied during training. Moreover, such networks may not be available for all target domains. To avoid overfitting into G_s and encourage the easier transfer of the optimized latent code to G_t , we use \mathcal{W} space optimization for this step. Finally, we initialize this w code for G_t and use additional attribute classifier loss [26] for T_t domain along with depth regularization $R(D)$ (Eq. 4). As an approximation, we assume that attribute classifier [25, 26] generalizes across all domains. We use $\mathcal{W}/\mathcal{W}+$ space optimization to control the quality and diversity of the outputs.

Let x be the source image, let pixel-wise *MSE* loss be represented as $L_{mse}(x, w, G) = MSE(x, G(w, M(\theta', \phi', c', r')))$, and let *LPIPS* loss be represented as $L_{lips}(x, w, G) = LPIPS(x, G(w, M(\theta', \phi', c', r')))$ where camera parameters are determined by Sec. 3.1. Let $L_d(w)$ be the depth regularizer (Eq. 4) and $A_t(x, w, G)$ be the attribute classifier loss. We define the algorithm of projection of a single source image into 3D avatars in Algorithm 1.

Algorithm 1: Projection of single image into 3D Avatar.

Input: source image $x \in \mathbb{R}^{n \times n \times 3}$; G_s, G_t , gradient-based optimizer F' and F'' .
Output: the embedded code w' for G_t

- 1 Initialize() the code $w = w_{avg}$;
- 2 **while not converged do**
- 3 $L \leftarrow L_{mse}(x, w, G_s) + L_{lips}(x, w, G_s) +$
 $L_d(w) + A_t(x, w, G_s)$;
- 4 $w \leftarrow w - \eta F'(\nabla_w L, w)$;
- 5 **end**
- 6 Initialize() the code $w' = w$;
- 7 **while not converged do**
- 8 $L' \leftarrow L_{mse}(x, w', G_t) + L_{lips}(x, w', G_t) +$
 $L_d(w') + A_t(x, w', G_t)$;
- 9 $w' \leftarrow w' - \zeta F''(\nabla_{w'} L', w')$;
- 10 **end**

Editing and Animation. Since our 3D domain adaptation is designed to preserve the properties of \mathcal{W} and \mathcal{S} spaces, we can perform semantic edits via InterFaceGAN [54], GANSpace [23], StyleSpace [65] *etc.*, and geometric edits using *TPS* (Sec. 3.4) and Δs interpolation (Sec. 3.2). To perform video editing, we design an encoder for EG3D based on *e4e* [61] to encode videos and transfer the edits from G_s to G_t based on the w codes [4, 6, 62]. We leave a more fine-grained approach for video processing as future work.



Figure 10. **Deformations using TPS.** Geometric edits using our proposed *TPS* (Thin Plate Spline) module learned on the frontal tri-plane features. Each sub-figure shows a 3D avatar and three examples of *TPS* deformations sampled from the learned 3D deformation space.

5. Results

5.1. Quantitative Results

In this section, we consider three important evaluations to verify the quality of the texture, geometry, and identity preservation in the new domain using the Caricature, Cartoons, and Pixar toons datasets. We also evaluate the ablation of our design choices and conduct a user study to assess the quality of generated avatars. In the evaluation, let G_{base} be the baseline naïve fine-tuning method which is trained with all the parameters using the losses in EG3D fine-tuned from FFHQ trained prior G_s . Note here we still align the cameras in G_{base} using the method defined in Sec. 3.1 and use adaptive discriminator [29] with R1 regularization for a fair comparison.

Texture Quality. To verify the quality of the texture, diversity of samples as well as to some extent, the geometry in the target domain T_t , we compare the *FID* [24] scores using G_{base} and G_t in Table 1. Note that in the case of Caricatures, we report two scores *i.e.* with and without using the attribute classifier loss in the training as discussed in Sec. 4. Notice that our method outperforms the naïve baseline method by a huge margin in some cases, especially in Caricatures and Cartoons. We attribute these differences to the mode collapse prone training of G_{base} which is correlated with flat geometry degenerate solution. We show visual results of the flat geometries learned by G_{base} and comparison in Fig. 2.

Geometric Quality. To quantify the flat geometries, in Table 2, we show three scores that help us understand such degenerate solutions. Here we consider coupled depth maps generated from sampling in the domains T_s (G_s) and T_t (G_t and G_{base}). First, we compute the expectation of the absolute mean differences (M_d) of the corresponding foreground depth maps sampled from T_s and T_t . We also compute the expectation of the absolute standard deviation differences (S_d) for the same setting. Here, we assume that the flatter geometries have a large difference in the depth maps as compared to the prior as indicated by M_d . Moreover, S_d computes the distance in the distribution of the depth values, where a larger difference indicates a narrow distribution, and hence a flatter geometry. We also notice

Table 1. **FID Computation.** FID (Fréchet Inception Distance) between the 2D dataset and the samples generated by the fine-tuned 3D GAN using baseline (G_{base}) and Ours (G_t). '*' represents the score with the inclusion of the attribute classifier loss discussed in Sec. 3.2.

Method	Caricatures	Cartoons	Pixar Toons
G_{base}	67.8	79.0	15.1
G_t (Ours)	19.4/20.2*	12.8	12.4

Table 2. **Geometry Evaluation.** Comparing the geometry using baseline method (G_{base}) and Ours (G_t). For the definition of M_d , S_d and $R(T_2)$, refer to Sec. 5.1.

Metric	Method	Caricatures	Cartoons	Pixar
$M_d \downarrow$	G_{base}	0.47	0.21	0.29
	G_t (Ours)	0.21	0.13	0.13
$S_d \downarrow$	G_{base}	0.22	0.14	0.15
	G_t (Ours)	0.15	0.10	0.09
$R(T_2) \downarrow$	G_{base}	2.99	3.39	4.01
	G_t (Ours)	2.27	1.62	1.56

that the flat geometry is correlated with the generator learning diverse poses when images are rendered under standard canonical camera parameters *i.e.* $M(0, 0, c, r)$. We hypothesize in the case of the flatter geometries, the model learns to pose information in the earlier layers instead of being camera view-dependent. To quantify this, since pose information may not be available for some domains (*e.g.* cartoons), we compute the $R(T_2)$ scores between corresponding images in the domain T_s (G_s) and T_t (G_t and G_{base}). Note that these scores are computed without the *TPS* module. Our scores are lower in all three metrics, hence, validating that our method avoids the degenerate solution and preserves the geometric distribution of the prior.

Identity Preservation. Identity preservation score is another important evaluation to check the quality of latent space linking between G_s and G_t . In Table 3, we compute the attribute loss (*BCE* loss) between the domains T_s and T_t

Table 3. **Identity Preservation.** Identity preservation using baseline (G_{base}) and Ours (G_t).

Method	Caricatures	Cartoons	Pixar Toons
G_{base}	1.28	0.92	0.85
G_t (Ours)	0.87	0.81	0.73



Figure 11. **Local edits.** Local edits performed on the 3D avatars using the S space.

using the attribute classifiers [25, 26]. Note that our method is able to preserve the identity better across the domains.

Ablation Study. In order to validate the importance of the losses, and components of our design choices, in Table 4, we show an ablation study of these components with regularizers. Note that we evaluate these design choices on the caricature dataset. Notice adding each component improves the corresponding scores in FID, M_d , S_d , and ID as discussed. Notice that by adding the TPS module, the FID

Table 4. **Ablation.** Ablation of the design choices made in Sec. 3. cam stands for the analysis in Sec. 3.1, Reg stands for the model after applying Eq. 2, DReg stands for the model after applying Eq. 4, and TPS stands for the model after applying Eq. 5 - 7. This is the G_t used in the comparison. Note that by adding TPS the scores are affected as the geometry is exaggerated *e.g.* the identity is affected. This module can be added to do geometry editing.

Method	FID	M_d	S_d	ID
$G_{\text{base}} - \text{cam}$	90.8	0.47	0.33	1.348
G_{base}	67.8	0.47	0.22	1.272
+ Reg	19.0	0.22	0.22	0.889
+ DReg	19.4	0.21	0.15	0.879
+ TPS	20.6	0.25	0.20	0.924

Table 5. **Ablation of TPS on metrics based on facial keypoints.**

Metric	with TPS	without TPS
Avg. Keypoint Distance	5.7	4.3
Avg. Keypoint Variation	0.06	0.04

is still comparable to G_t . The slight drop is attributed to the stretching and squeezing of some parts of the texture (See Fig. 10). Nevertheless, by adding this module, we achieve better control over geometry and produce exaggerated features for a small drop in texture quality. In order to show that the TPS transformations are not random, we compute the FID scores with randomly perturbed front plane features which are derived from the perturbations of control points near the face *e.g.* taken at the early stages of training after it has stabilized a bit. This setup has less perturbation in the background. We found that the FID score is worse *i.e.* 25.5 hence validating non-random transformations.

The purpose of adding the TPS module is to model the exaggerated geometries in the data and at the same time achieve geometric editing and animation capabilities (see Project Page). We ablate (see Table 5) the average face Keypoint Distance between the paired FFHQ-generated images and corresponding avatars and the face Keypoint Variation (standard deviation of keypoints) with and without the TPS module. The results show that with TPS module, the deviations are large and match the exaggerations.

User Study. For results on real images refer to Fig. 7. We conducted a user study using 50 real images (25 caricatures, and 25 Pixar) on identity preservation and 3D consistency versus the baseline method. We asked 21 unique workers where each triplet was reviewed by 10 workers and our avatars were chosen **92%** of the time.

5.2. Qualitative Results

For qualitative results, we show the results of the domain adaptation, as well as the personalized edits (geometric and semantic), performed on the resultant 3D avatars. First, in

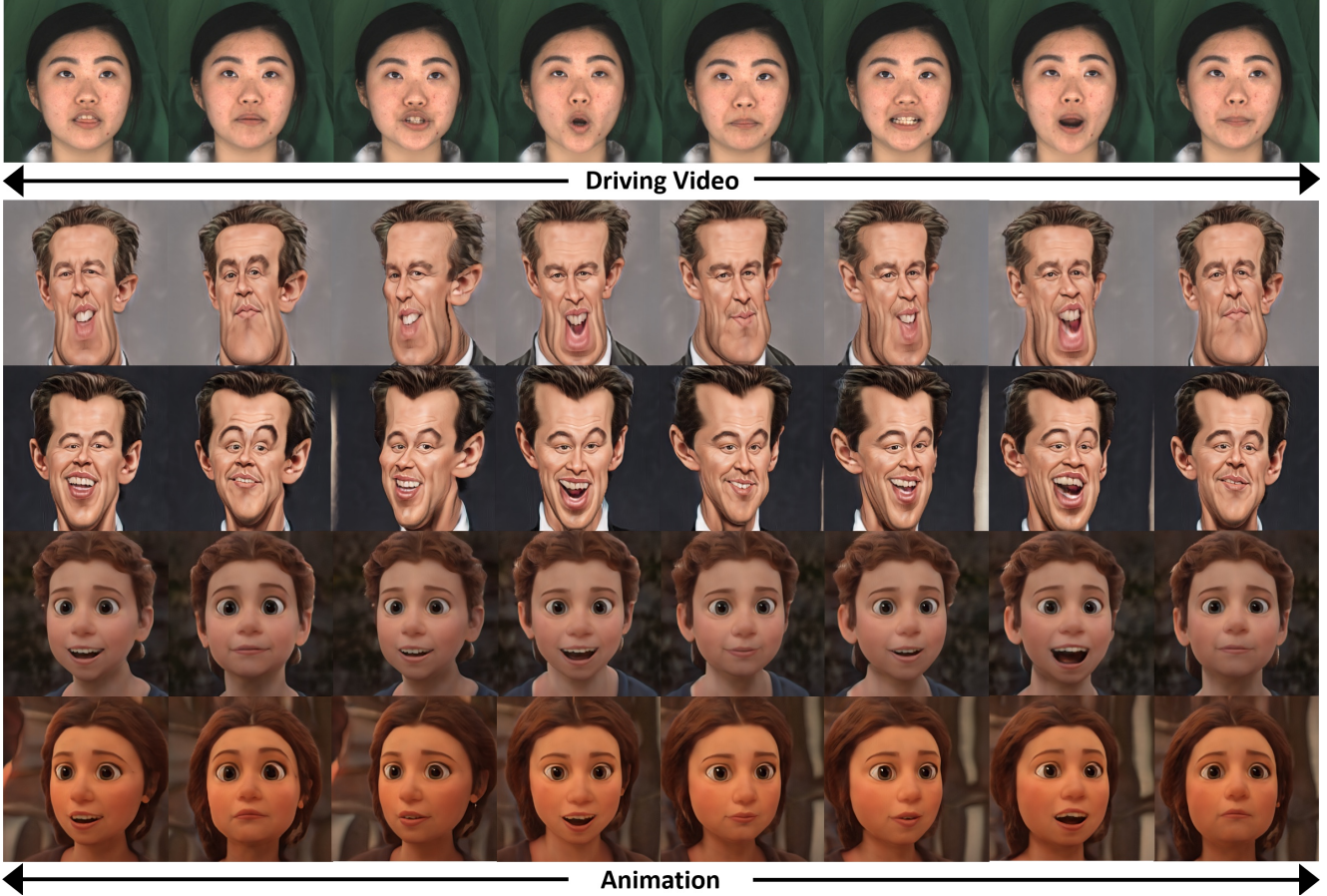


Figure 12. **3D avatar animation.** Animation of 3D avatars generated using a driving video encoded in source domain T_s and applied to samples in target domain T_t . The top row shows the driving video and the subsequent rows show generated animations using a random Caricature or Pixar toon. The head pose is changed in each frame of the generated animation to show 3D consistency.

order to show the quality of domain adaptation, identity preservation, and geometric consistency, in Fig. 6, we show results from G_s and corresponding results from 3D avatar generator G_t trained on Caricatures, Pixar toons, Cartoons, and Comic domains. Next, in order to show that the method generalizes to real images, we use the method described in Sec. 4 to project and transfer the latent code from G_s to G_t to produce the 3D avatars. In Fig. 8, we show our results of real to 3D avatar transfer. Notice the quality both in terms of texture as well as geometry for both these results achieved by our method. Next, we show geometric and semantic edits possible to produce personalized 3D avatars:

Geometry Edits. We show two type of geometric edits *i.e.* Δs interpolation (Sec. 3.2) and deformation using *TPS* (Sec. 3.4). First, in Fig. 9, we show the geometry interpolation by interpolating between original s activations of G_s and learned Δs parameters. In Fig. 10, we show some additional exaggerations in caricatures using the learned 3D deformation latent space of *TPS* module.

Semantic Edits and Animation. Since in our method, we

encourage the latent regularization to preserve the properties of the latent space learned by the G_s generator, in Fig. 11 we show S space edits performed on the 3D avatars. Notice the quality of edits in terms of locality and adaptability. Additionally, we can edit semantics like hair as opposed to 3D morphable model based methods. In Fig. 12, thanks to the latent space semantics preservation ensured by our method, we can perform some video edits to create a coherent animation based on the difference of w codes of video encoded in G_s (Sec. 4) and applied to layers 7 – 10 in G_t . Notice the quality of expressions, identity preservation, and 3D consistency across each identity in each row.

6. Conclusion

We tackled two open research problems in this paper. In the first part, we proposed the first domain adaptation method for 3D-GANs to the best of our knowledge. This part yields two linked EG3D generators, one in the photo-realistic source domain of faces, and another EG3D generator in an artistic target domain. As possible target domains,

we show results for cartoons, caricatures, and comics. In the second part, we built on domain adaptation to create 3D avatars in an artistic domain that can be edited and animated. Our framework consists of multiple technical components introduced in this paper. First, we propose a technique for camera space estimation for artistic domains. Second, we introduce a set of regularizers and loss functions that can regularize the fine-tuning of EG3D in such a way that enough of the 3D structure and geometry of the original model is kept, while the distinguishing attributes of the artistic domain, such as textures and colors and local geometric deformations can still be learned. Third, we introduce a geometric deformation module that can reintroduce larger geometric deformations in a controlled manner. These larger geometric deformations can interact and cooperate with EG3D so that semantic edits are still possible. Finally, we propose an embedding algorithm that is especially suitable for two linked EG3D generator networks.

7. Limitations

Our method also has some limitations. Overall, the visual quality is limited by the quality of StyleGAN2 pretraining. While we found the quality to be very high for the datasets shown in the paper, it relies on hundreds of images in the target domain to be available. Would be interesting to do few-shot domain adaptation in the future. Further, edits are largely limited to semantic edits of EG3D and global space deformations by *TPS*. Our method does not enable fine-grained geometric edits. Finally, a large part of our method is face specific. We justify this specialization by the importance of human models and the specific target domain of editable 3D avatars. We nevertheless believe that domain adaptation of general 3D-GANs will be an interesting avenue of future work.

8. Ethical Concerns

Deep learning-based image and video processing is a tool for image/video understanding, animation, and artistic expression. Similar to most software in this domain, our work could be used to produce offensive results. An application of concern would be if a user would generate offensive caricatures and cartoons of other people without consent. This can be used to insinuate biases against people or can have a detrimental effect on a person’s autonomy, dignity, and/or privacy. The images used in this work are taken or derived from the FFHQ [42] dataset which has appropriate licenses for non-commercial research use and to the best of our knowledge, is committed to protecting the privacy of the individuals who do not wish to be included.

9. Training Details

We train our models on 4 V100 GPUs with a batch size of 8. Similar to EG3D, we start training from the neural rendering resolution of 64^2 which is increased during the training. Then we do fine-tuning on 128^2 resolution to produce the final 512^2 outputs. We sample 100k samples from each dataset. We train Caricatures on ~ 880 *kimgs*, Pixar, and Cartoons on ~ 500 *kimgs*. We fine-tune these models on 128^2 neural rendering resolution for an additional 80 – 160 *kimgs*. Other hyperparameters of learning rates are the same as the EG3D. We train the *TPS* module on ~ 2000 *kimgs*. We set the weight for the regularization term $R(\Delta s)$ as 0.001, and $R(D)$ as 0.005. For the *TPS* training we use the weights for α , β , σ and $R(T_2)$ as 150, 1, 3 and 1 respectively. For inversion, we perform 200 steps for the source domain inversion and 400 steps for the target domain to generate the final avatar.

Table 6. **FID comparison with SCG: StyleCariGAN, DSG: DualStyleGAN.**

Method	Cari. (SCG)	Cari. (DSG)	Pixar (DSG)	Cartoons (DSG)
2D-GANs	51.68	96.49	166.78	105.57
Ours	62.90	89.73	162.06	111.35

10. Comparison to 2D-GANs

The quality drop is expected when the final model is compared with 2D-GANs as we derive our datasets from these 2D-GANs fine-tuned on avatar datasets. In Table 6, we compute the FID scores between the datasets (size ~ 200 images for DualStyleGAN dataset) used to train these 2D-GANs and our corresponding avatar generators (size $\sim 10k$ images). We used the *3DCaricShop* [48] dataset for Caricatures as the authors of *WebCaricature* did not reply with the download link. The scores are comparable, even better in the case of Caricatures and Pixar, probably due to our regularizers including 3D view consistency.

11. Importance of Front Tri-plane Features

As discussed in Sec. 3.4 of the main paper, the front tri-plane of the EG3D architecture encodes most of the texture and depth information in the output. In Fig. 13, we show two images with their front and other side plane swapped. Then we show the corresponding effect on the output image. Notice that the results are consistent with the analysis in Sec. 3.4 where the front tri-plane dominates the information for output texture and depth.

12. Stylization

To validate that our chosen layers in Sec. 3.2 are responsible for geometrical and texture changes, we resort to a stylization technique. For stylization given an arbitrary reference image *e.g.* painting, we use the *Style Loss* [19] to update the layers of G_s or G_t . Essentially, we use the same parameters used in Sec. 3.1. A critical technique to achieve multi-view consistency and circumvent the ghosting face artifact due to single image overfitting is to rotate the camera to cover the θ' and ϕ' ranges in Sec. 3.1 in the main paper uniformly during the optimization. In Fig. 14, we show some results using only the layers used in *Texture Regularization* (Sec. 3.2). Note the high-quality texture change in the images. In Fig. 15, we show results by adding layers of *Geometry Regularization* (Sec. 3.2). Note that the geometry is changed in the examples when we use this module. Note that in this example, the geometry is not expected to match as there is no such loss in the optimization. This example results in some arbitrary geometry change that is not flat. This validates our choice of geometry and texture layers used in this paper.

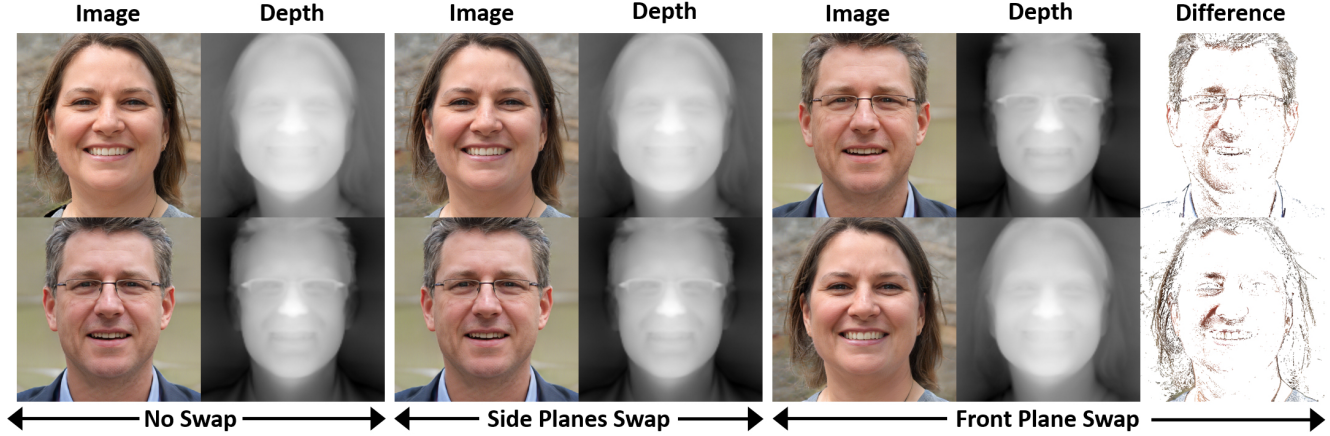


Figure 13. **Swapping tri-planes.** Validation of the information stored in the front tri-plane. Given two images and their tri-plane representations, there is almost no change in the final output if the side planes are swapped. While the output changes completely when the front tri-plane is swapped.

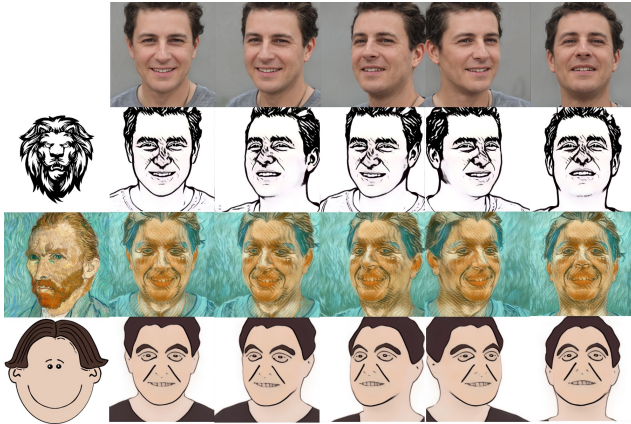


Figure 14. **Validation of texture regularization.** Style Transfer using the texture layers discussed in the main paper. Here the style of the image is changed using the texture layers.

13. Failure cases

Our method has some failure cases that stem from the samples of the 2D-GANs *i.e.* DualStyleGAN [68] and StyleCariGAN [26]. In the caricature domain, the random samples generated in the case of G_t trained on StyleCariGAN samples can have some severe artifacts (See Fig. 16). Although these do not appear often, such a sample can be improved by attribute classifier and depth regularization losses used on a single image as discussed in Sec. 4 of the main paper.

14. Depth Map visualization

In order to show some more samples and the corresponding depth maps for G_{base} and G_t , in Fig. 17, we show some



Figure 15. **Validation of geometry regularization.** Geometry change (third row) using the geometry layers discussed in the main paper. Here the style (second row) is changed using the texture layers and the geometry (third row) is changed using the geometry layers. Note that the geometry is not correct as we apply *Style Loss* using a single image and is only used to demonstrate the usage of different layers.



Figure 16. **Failure cases.** Failure cases in caricature generation that stems from the artifacts in the 2D-GAN from which the dataset is generated. Here the artifacts are the result of the generated results of StyleCariGAN [26]

samples from both the generators and the corresponding depth maps with pose changes. Notice the flat geometry in the case of G_{base} results. Next, in Fig. 18 and Fig. 19, we show some grid samples of our method with depth maps on the Caricature, Pixar toon, Cartoon, and Comic datasets.

15. Video Results

We also show our 3D avatar editing results in videos. We design a simple UI to show that the avatars can be edited in an interactive manner. Please refer to the [webpage](#) for editing videos and the interactive editing sessions.



Figure 17. **Comparison with Naive method.** Results of the Caricatures and Pixar toons were viewed from different viewpoints and compared with the baseline method. Note that the depth maps are also visualized highlighting flat geometry. For more results refer to the videos on the [Project Page](#).



Figure 18. **Grid samples.** Samples from the source domain and corresponding results in the target domain. Corresponding images and depth outputs of the Caricatures and Pixar Toons are shown.



Figure 19. **Grid samples.** Extension of Fig. 18. Corresponding images and depth outputs of the Comics and Cartoons are shown.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, Seoul, Korea, 2019. IEEE. 2, 7
- [2] Rameen Abdal, Peihao Zhu, John Femiani, Niloy Mitra, and Peter Wonka. Clip2stylegan: Unsupervised extraction of stylegan edit directions. In *ACM SIGGRAPH 2022 Conference Proceedings*, SIGGRAPH '22, New York, NY, USA, 2022. Association for Computing Machinery. 2
- [3] Rameen Abdal, Peihao Zhu, Niloy J. Mitra, and Peter Wonka. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *ACM Trans. Graph.*, 40(3), may 2021. 2, 3
- [4] Rameen Abdal, Peihao Zhu, Niloy J. Mitra, and Peter Wonka. Video2stylegan: Disentangling local and global variations in a video, 2022. 7
- [5] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder via iterative refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021. 2
- [6] Yuval Alaluf, Or Patashnik, Zongze Wu, Asif Zamir, Eli Shechtman, Dani Lischinski, and Daniel Cohen-Or. Third time's the charm? image and video editing with stylegan3. *CoRR*, abs/2201.13433, 2022. 7
- [7] Yuval Alaluf, Omer Tov, Ron Mokady, Rinon Gal, and Amit H. Bermano. Hyperstyle: Stylegan inversion with hypernetworks for real image editing. *CoRR*, abs/2111.15666, 2021. 2
- [8] Anonymous. 3d generation on imagenet. In *Open Review*, 2023. 3
- [9] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 2
- [10] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *arXiv*, 2021. 1, 2, 3
- [11] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks, 2021. 2, 4
- [12] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5799–5809, 2021. 2
- [13] Yen-Chi Cheng, Chieh Hubert Lin, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, and Ming-Hsuan Yang. Inout: Diverse image outpainting via gan inversion. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [14] Yunje Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [15] Min Jin Chong and David A. Forsyth. Jojogan: One shot face stylization. *CoRR*, abs/2112.11641, 2021. 2
- [16] Min Jin Chong, Hsin-Ying Lee, and David Forsyth. Stylegan of all trades: Image manipulation with only pretrained stylegan. *arXiv preprint arXiv:2111.01619*, 2021. 2
- [17] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *arXiv preprint arXiv:2108.00946*, 2021. 2
- [18] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators, 2021. 2
- [19] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016. 12
- [20] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. 2
- [21] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *International Conference on Learning Representations*, 2022. 2
- [22] Fangzhou Han, Shuquan Ye, Mingming He, Menglei Chai, and Jing Liao. Exemplar-based 3d portrait stylization. *IEEE Transactions on Visualization and Computer Graphics*, 2021. 2
- [23] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. *arXiv preprint arXiv:2004.02546*, 2020. 2, 7
- [24] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 8
- [25] Jing Huo, Wenbin Li, Yinghuan Shi, Yang Gao, and Hujun Yin. Webcaricature: a benchmark for caricature recognition. In *British Machine Vision Conference*, 2018. 3, 7, 9
- [26] Wonjong Jang, Gwangjin Ju, Yucheol Jung, Jiaolong Yang, Xin Tong, and Seungyong Lee. Stylecarigan: Caricature generation via stylegan feature map modulation. 40(4), 2021. 2, 3, 7, 9, 13
- [27] Yucheol Jung, Wonjong Jang, Soongjin Kim, Jiaolong Yang, Xin Tong, and Seungyong Lee. Deep deformable 3d caricatures with learned shape control. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings*. ACM, aug 2022. 2
- [28] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation, 2017. 2
- [29] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Proc. NeurIPS*, 2020. 1, 2, 6, 8

- [30] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *arXiv preprint arXiv:2006.06676*, 2020. 3
- [31] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks, 2021. 1, 2
- [32] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1, 3
- [33] Tero Karras, Samuli Laine, and Timo Aila. A Style-Based generator architecture for generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 43(12):4217–4228, Dec. 2021. 2
- [34] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020. 2
- [35] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 5, 7
- [36] Thomas Leimkühler and George Drettakis. Freestylegan: Free-view editable portrait rendering with the camera manifold. 40(6), 2021. 3
- [37] Chieh Hubert Lin, Hsin-Ying Lee, Yen-Chi Cheng, Sergey Tulyakov, and Ming-Hsuan Yang. Infinitygan: Towards infinite-pixel image synthesis. In *International Conference on Learning Representations (ICLR)*, 2022. 2
- [38] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 2
- [39] Sangwoo Mo, Minsu Cho, and Jinwoo Shin. Freeze the discriminator: a simple baseline for fine-tuning gans, 2020. 2
- [40] Michael Niemeyer and Andreas Geiger. Campari: Camera-aware decomposed generative neural radiance fields. In *2021 International Conference on 3D Vision (3DV)*, pages 951–961. IEEE, 2021. 2
- [41] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 2
- [42] NVlabs. ffhq-dataset. <https://github.com/NVlabs/ffhq-dataset>. 12
- [43] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. StyleSDF: High-Resolution 3D-Consistent Image and Geometry Generation. *arXiv preprint arXiv:2112.11427*, 2021. 1, 2
- [44] Xingang Pan, Bo Dai, Ziwei Liu, Chen Change Loy, and Ping Luo. Do 2d gans know 3d shape? unsupervised 3d shape reconstruction from 2d image gans. *arXiv preprint arXiv:2011.00844*, 2020. 2
- [45] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 2
- [46] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery, 2021. 2
- [47] Justin N. M. Pinkney and Doron Adler. Resolution dependent gan interpolation for controllable image synthesis between domains, 2020. 2, 3
- [48] Yuda Qiu, Xiaojie Xu, Lingteng Qiu, Yan Pan, Yushuang Wu, Weikai Chen, and Xiaoguang Han. 3dcaricshop: A dataset and a baseline method for single-view 3d caricature face reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10236–10245, 2021. 12
- [49] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021. 2
- [50] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks, 2015. 2
- [51] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. *arXiv preprint arXiv:2008.00951*, 2020. 2
- [52] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *arXiv preprint arXiv:2106.05744*, 2021. 2
- [53] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [54] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2, 7
- [55] Yichun Shi, Divyansh Aggarwal, and Anil K Jain. Lifting 2d stylegan for 3d-aware face generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6258–6266, 2021. 2
- [56] Ivan Skorokhodov, Aliaksandr Siarohin, Yinghao Xu, Jian Ren, Hsin-Ying Lee, Peter Wonka, and Sergey Tulyakov. 3d generation on imagenet. In *International Conference on Learning Representations (ICLR)*, 2023. 2
- [57] Guoxian Song, Linjie Luo, Jing Liu, Wan-Chun Ma, Chunpong Lai, Chuanxia Zheng, and Tat-Jen Cham. Agilegan: Stylizing portraits by inversion-consistent transfer learning. *ACM Trans. Graph.*, 40(4), jul 2021. 2
- [58] Jingxiang Sun, Xuan Wang, Yichun Shi, Lizhen Wang, Jue Wang, and Yebin Liu. Ide-3d: Interactive disentangled editing for high-resolution 3d-aware portrait synthesis, 2022. 1, 2

- [59] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt. Stylerig: Rigging stylegan for 3d control over portrait images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6142–6151, 2020. 2
- [60] Ayush Tewari, Mohamed Elgharib, Mallikarjun BR, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt. Pie: Portrait image embedding for semantic control. volume 39, December 2020. 2
- [61] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *arXiv preprint arXiv:2102.02766*, 2021. 2, 7
- [62] Rotem Tzaban, Ron Mokady, Rinon Gal, Amit H. Bermano, and Daniel Cohen-Or. Stitch it in time: Gan-based facial editing of real videos. *CoRR*, abs/2201.08361, 2022. 7
- [63] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Cross-domain and disentangled face manipulation with 3d guidance. *IEEE Transactions on Visualization and Computer Graphics*, 2022. 2
- [64] WarBean. tps-stn-pytorch. https://github.com/WarBean/tps_stn_pytorch. 6
- [65] Zongze Wu, Dani Lischinski, and Eli Shechtman. Stylespace analysis: Disentangled controls for stylegan image generation. *arXiv preprint arXiv:2011.12799*, 2020. 2, 5, 7
- [66] Yinghao Xu, Menglei Chai, Zifan Shi, Sida Peng, Ivan Skorokhodov, Aliaksandr Siarohin, Ceyuan Yang, Yujun Shen, Hsin-Ying Lee, Bolei Zhou, et al. Discoscene: Spatially disentangled generative radiance fields for controllable 3d-aware scene synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 2
- [67] Yinghao Xu, Sida Peng, Ceyuan Yang, Yujun Shen, and Bolei Zhou. 3d-aware image synthesis via learning structural and textural representations. *arXiv preprint arXiv:2112.10759*, 2021. 1, 2
- [68] Shuai Yang, Liming Jiang, Ziwei Liu, and Chen Change Loy. Pastiche master: Exemplar-based high-resolution portrait style transfer. In *CVPR*, 2022. 2, 3, 13
- [69] Zipeng Ye, Mengfei Xia, Yanan Sun, Ran Yi, Minjing Yu, Juyong Zhang, Yu-Kun Lai, and Yong-Jin Liu. 3d-CariGAN: An end-to-end solution to 3d caricature generation from normal face photos. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2021. 2
- [70] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 2
- [71] Kai Zhang, Gernot Riegler, Noah Snaveley, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [72] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image editing. In *European Conference on Computer Vision*, pages 592–608. Springer, 2020. 2
- [73] Peihao Zhu, Rameen Abdal, John Femiani, and Peter Wonka. Mind the gap: Domain gap control for single shot domain adaptation for generative adversarial networks. In *International Conference on Learning Representations*, 2022. 2
- [74] Peihao Zhu, Rameen Abdal, Yipeng Qin, John Femiani, and Peter Wonka. Improved stylegan embedding: Where are the good latents?, 2020. 2
- [75] zllrunning. face-parsing.pytorch. <https://github.com/zllrunning/face-parsing.PyTorch>. 2, 3