

3D-aware Conditional Image Synthesis

Kangle Deng Gengshan Yang Deva Ramanan Jun-Yan Zhu
Carnegie Mellon University

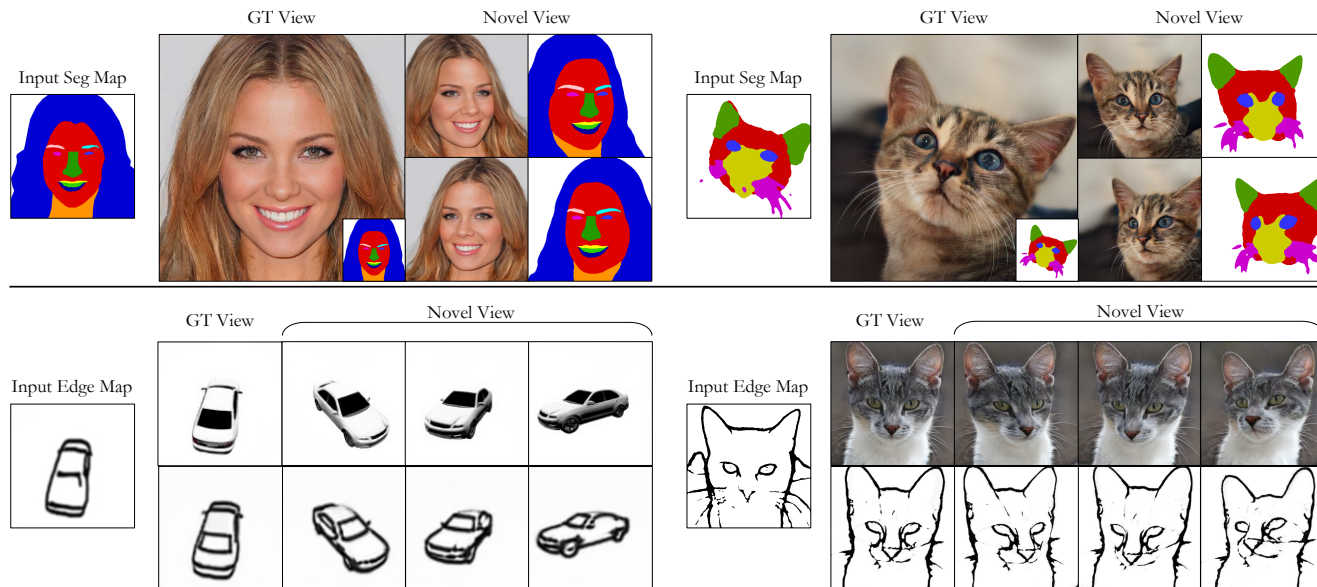


Figure 1. Given a 2D label map as input, such as a segmentation or edge map, our model learns to predict high-quality 3D labels, geometry, and appearance, which enables us to render both labels and RGB images from different viewpoints. The inferred 3D labels further allow interactive editing of label maps from any viewpoint, as shown in Figure 10.

Abstract

We propose *pix2pix3D*, a 3D-aware conditional generative model for controllable photorealistic image synthesis. Given a 2D label map, such as a segmentation or edge map, our model learns to synthesize a corresponding image from different viewpoints. To enable explicit 3D user control, we extend conditional generative models with neural radiance fields. Given widely-available posed monocular image and label map pairs, our model learns to assign a label to every 3D point in addition to color and density, which enables it to render the image and pixel-aligned label map simultaneously. Finally, we build an interactive system that allows users to edit the label map from different viewpoints and generate outputs accordingly.

1. Introduction

Content creation with generative models has witnessed tremendous progress in recent years, enabling high-quality,

user-controllable image and video synthesis [19, 20, 24, 34]. In particular, image-to-image translation methods [29, 56, 86] allow users to interactively create and manipulate a high-resolution image given a 2D input label map. Unfortunately, existing image-to-image translation methods operate purely in 2D, without explicit reasoning of the underlying 3D structure of the content. As shown in Figure 1, we aim to make conditional image synthesis 3D-aware, allowing not only 3D content generation but also viewpoint manipulation and attribute editing (e.g., car shape) in 3D.

Synthesizing 3D content conditioned on user input is challenging. For model training, it is costly to obtain large-scale datasets with paired user inputs and their desired 3D outputs. During test time, 3D content creation often requires multi-view user inputs, as a user may want to specify the details of 3D objects using 2D interfaces from different viewpoints. However, these inputs may not be 3D-consistent, providing conflicting signals for 3D content creation.

To address the above challenges, we extend conditional generative models with 3D neural scene representations. To

enable *cross-view* editing, we additionally encode semantic information in 3D, which can then be rendered as 2D label maps from different viewpoints. We learn the aforementioned 3D representation using only 2D supervision in the form of image reconstruction and adversarial losses. While the reconstruction loss ensures the alignment between 2D user inputs and corresponding 3D content, our pixel-aligned conditional discriminator encourages the appearance and labels to look plausible while remaining pixel-aligned when rendered into novel viewpoints. We also propose a cross-view consistency loss to enforce the latent codes to be consistent from different viewpoints.

We focus on 3D-aware semantic image synthesis on the CelebAMask-HQ [38], AFHQ-cat [16], and shapenet-car [10] datasets. Our method works well for various 2D user inputs, including segmentation maps and edge maps. Our method outperforms several 2D and 3D baselines, such as Pix2NeRF variants [6], SofGAN [11], and SEAN [89]. We further ablate the impact of various design choices and demonstrate applications of our method, such as cross-view editing and explicit user control over semantics and style. Please see our [website](#) for more results and [code](#).

2. Related Work

Neural Implicit Representation. Neural implicit fields, such as DeepSDF and NeRFs [46, 54], model the appearance of objects and scenes with an implicitly defined, continuous 3D representation parameterized by neural networks. They have produced significant results for 3D reconstruction [67, 90] and novel view synthesis applications [39, 43, 44, 48, 81] thanks to their compactness and expressiveness. NeRF and its descendants aim to optimize a network for an individual scene, given hundreds of images from multiple viewpoints. Recent works further reduce the number of training views through learning network initializations [13, 70, 79], leveraging auxiliary supervision [18, 30], or imposing regularization terms [50]. Recently, explicit or hybrid representations of radiance fields [12, 48, 61] have also shown promising results regarding quality and speed. In our work, we use hybrid representations for modeling both user inputs and outputs in 3D, focusing on synthesizing novel images rather than reconstructing an existing scene. A recent work Pix2NeRF [6] aims to translate a single image to a neural radiance field, which allows single-image novel view synthesis. In contrast, we focus on 3D-aware user-controlled content generation.

Conditional GANs. Generative adversarial networks (GANs) learn the distribution of natural images by forcing the generated and real images to be indistinguishable. They have demonstrated high-quality results on 2D image synthesis and manipulation [1, 3, 5, 20, 33–35, 59, 65, 72, 84, 85]. Several methods adopt image-conditional GANs [29, 47] for

user-guided image synthesis and editing applications [26, 27, 38, 40, 55, 56, 62, 74, 86, 89]. In contrast, we propose a 3D-aware generative model conditioned on 2D user inputs that can render view-consistent images and enable interactive 3D editing. Recently, SoFGAN [11] uses a 3D semantic map generator and a 2D semantic-to-image generator to enable 3D-aware generation, but using 2D generators does not ensure 3D consistency.

3D-aware Image Synthesis. Early data-driven 3D image editing systems can achieve various 3D effects but often require a huge amount of manual effort [14, 37]. Recent works have integrated the 3D structure into learning-based image generation pipelines using various geometric representations, including voxels [22, 88], voxelized 3D features [49], and 3D morphable models [71, 78]. However, many rely on external 3D data [71, 78, 88]. Recently, neural scene representations have been integrated into GANs to enable 3D-aware image synthesis [8, 9, 21, 51–53, 64, 77]. Intriguingly, these 3D-aware GANs can learn 3D structures without any 3D supervision. For example, StyleNeRF [21] and EG3D [8] learn to generate 3D representations by modulating either NeRFs or explicit representations with latent style vectors. This allows them to render high-resolution view-consistent images. Unlike the above methods, we focus on conditional synthesis and interactive editing rather than random sampling. Several works [17, 28, 42, 76] have explored sketch-based shape generation but they do not allow realistic image synthesis.

Closely related to our work, Huang et al. [25] propose synthesizing novel views conditional on a semantic map. Our work differs in three ways. First, we can predict full 3D labels, geometry, and appearance, rather than only 2D views, which enables cross-view editing. Second, our method can synthesize images with a much wider baseline than Huang et al. [25]. Finally, our learning algorithm does not require ground truth multi-view images of the same scene. Two recent works, FENeRF [69] and 3DSGAN [80], also leverage semantic labels for training 3D-aware GANs, but they do not support conditional inputs and require additional efforts (e.g., GAN-inversion) to allow user editing. Three concurrent works, IDE-3D [68], NeRFFaceEditing [31], and sem2nerf [15], also explore the task of 3D-aware generation based on segmentation masks. However, IDE-3D and sem2nerf only allow editing on a fixed view, and NeRF-FaceEditing focuses on real image editing rather than generation. All of them do not include results for other input modalities. In contrast, we present a general-purpose method that works well for diverse datasets and input controls.

3. Method

Given a 2D label map I_s , such as a segmentation or edge map, `pix2pix3D` generates a 3D-volumetric representation of geometry, appearance, and labels that can be rendered from different viewpoints. Figure 2 provides an overview.

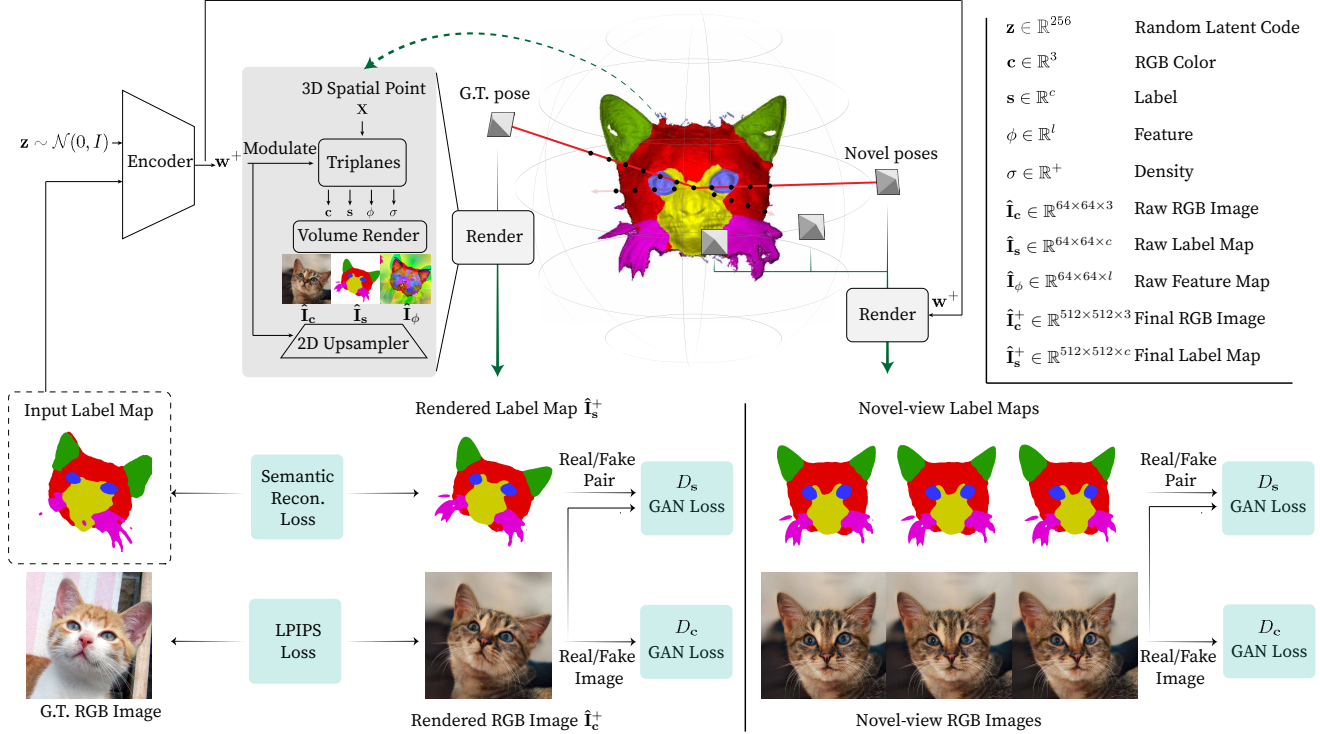


Figure 2. **Overall pipeline.** Given a 2D label map (e.g., segmentation map), a random latent code z , and a camera pose \hat{P} as inputs, our generator renders the label map and image from viewpoint \hat{P} . Intuitively, the input label map specifies the geometric structure, while the latent code captures the appearance, such as hair color. We begin with an encoder that encodes both the input label map and the latent code into style vectors w^+ . We then use w^+ to modulate our 3D representation, which takes a spatial point x and outputs (1) color $c \in \mathbb{R}^3$, (2) density σ , (3) feature $\phi \in \mathbb{R}^l$, and (4) label $s \in \mathbb{R}^c$. We then perform volumetric rendering and 2D upsampling to get the high-res label map \hat{I}_s^+ and RGB Image \hat{I}_c^+ . For those rendered from ground-truth poses, we compare them to ground-truth labels and images with an LPIPS loss and label reconstruction loss. We apply a GAN loss on labels and images rendered from both novel and original viewpoints.

We first introduce the formulation of our 3D conditional generative model for 3D-aware image synthesis in Section 3.1. Then, in Section 3.2, we discuss how to learn the model from color and label map pairs $\{I_c, I_s\}$ associated with poses P .

3.1. Conditional 3D Generative Models

Similar to EG3D [8], we adopt a hybrid representation for the density and appearance of a scene and use style vectors to modulate the 3D generations. To condition the 3D representations on 2D label map inputs, we introduce a conditional encoder that maps a 2D label map into a latent style vector. Additionally, pix2pix3D produces 3D labels that can be rendered from different viewpoints, allowing for cross-view user editing.

Conditional Encoder. Given a 2D label map input I_s and a random latent code sampled from the spherical Gaussian space $z \sim \mathcal{N}(0, I)$, our conditional encoder E outputs a list of style vectors $w^+ \in \mathbb{R}^{l \times 256}$,

$$w^+ = E(I_s, z),$$

where $l = 13$ is the number of layers to be modulated. Specifically, we encode I_s into the first 7 style vectors that

represent the global geometric information of the scene. We then feed the random latent code z through a Multi-Layer Perceptron (MLP) mapping network to obtain the rest of the style vectors that control the appearance.

Conditional 3D Representation. Our 3D representation is parameterized by tri-planes followed by a 2-layer MLP f [8], which takes in a spatial point $x \in \mathbb{R}^3$ and returns 4 types of outputs: (1) color $c \in \mathbb{R}^3$, (2) density $\sigma \in \mathbb{R}^+$, (3) feature $\phi \in \mathbb{R}^{64}$ for the purpose of 2D upsampling, and most notably, (4) label $s \in \mathbb{R}^c$, where c is the number of classes if I_s is a segmentation map, otherwise 1 for edge labels. We make the field conditional by modulating the generation of tri-planes F^{tri} with the style vectors w^+ . We also remove the view dependence of the color following [8, 21]. Formally,

$$(c, s, \sigma, \phi) = f(F_{w^+}^{\text{tri}}(x)).$$

Volume Rendering and Upsampling. We apply volumetric rendering to synthesize color images [32, 46]. In addition, we render label maps, which are crucial for enabling cross-view editing (Section 4.3) and improving rendering quality (Table 1). Given a viewpoint \hat{P} looking at the scene origin,

we sample N points along the ray that emanates from a pixel location and query density, color, labels, and feature information from our 3D representation. Let \mathbf{x}_i be the i -th sampled point along the ray r . Let \mathbf{c}_i , \mathbf{s}_i and ϕ_i be the color, labels, and the features of \mathbf{x}_i . Similar to [69], The color, label map, and feature images are computed as the weighted combination of queried values,

$$\hat{\mathbf{I}}_c(r) = \sum_{i=1}^N \tau_i \mathbf{c}_i, \quad \hat{\mathbf{I}}_s(r) = \sum_{i=1}^N \tau_i \mathbf{s}_i, \quad \hat{\mathbf{I}}_\phi(r) = \sum_{i=1}^N \tau_i \phi_i, \quad (1)$$

where the transmittance τ_i is computed as the probability of a photon traversing between the camera center and the i -th point given the length of the i -th interval δ_i ,

$$\tau_i = \prod_{j=1}^i \exp(-\sigma_j \delta_j) (1 - \exp(-\sigma_i \delta_i)).$$

Similar to prior works [8, 21, 52], we approximate Equation 1 by 2D Upsampler U to reduce the computational cost. We render high-res 512×512 images in two passes. In the first pass, we render low-res 64×64 images $\hat{\mathbf{I}}_c, \hat{\mathbf{I}}_s, \hat{\mathbf{I}}_\phi$. Then a CNN up-sampler U is applied to obtain high-res images,

$$\hat{\mathbf{I}}_c^+ = U(\hat{\mathbf{I}}_c, \hat{\mathbf{I}}_\phi), \quad \hat{\mathbf{I}}_s^+ = U(\hat{\mathbf{I}}_s, \hat{\mathbf{I}}_\phi).$$

3.2. Learning Objective

Learning conditional 3D representations from monocular images is challenging due to its under-constrained nature. Given training data of associated images, label maps, and camera poses predicted by an off-the-shelf model, we carefully construct learning objectives, including reconstruction, adversarial, and cross-view consistency losses. These objectives will be described below.

Reconstruction Loss. Given a ground-truth viewpoint \mathbf{P} associated with the color and label maps $\{\mathbf{I}_c, \mathbf{I}_s\}$, we render color and label maps from \mathbf{P} and compute reconstruction losses for both high-res and low-res output. We use LPIPS [82] to compute the image reconstruction loss \mathcal{L}_c for color images. For label reconstruction loss \mathcal{L}_s , we use the balanced cross-entropy loss for segmentation maps or L2 Loss for edge maps,

$$\mathcal{L}_{\text{recon}} = \lambda_c \mathcal{L}_c(\mathbf{I}_c, \{\hat{\mathbf{I}}_c, \hat{\mathbf{I}}_c^+\}) + \lambda_s \mathcal{L}_s(\mathbf{I}_s, \{\hat{\mathbf{I}}_s, \hat{\mathbf{I}}_s^+\}),$$

where λ_c and λ_s balance two terms.

Pixel-aligned Conditional Discriminator. The reconstruction loss alone fails to synthesize detailed results from novel viewpoints. Therefore, we use an adversarial loss [20] to enforce renderings to look realistic from random viewpoints. Specifically, we have two discriminators D_c and D_s for RGB images and label maps, respectively. D_c is a widely-used GAN loss that takes real and fake images as input, while

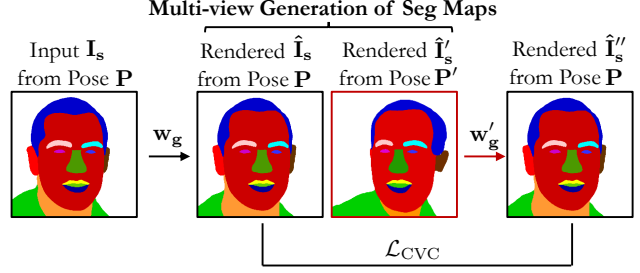


Figure 3. **Cross-View Consistency Loss.** Given an input label map \mathbf{I}_s and its associated pose \mathbf{P} , we first infer the geometry latent code w_g . From w_g , we can generate a label map $\hat{\mathbf{I}}_s$ from the same pose \mathbf{P} , and $\hat{\mathbf{I}}'_s$ from a random pose \mathbf{P}' . Next, we infer w'_g from the novel view $\hat{\mathbf{I}}'_s$, and render it back to the original pose \mathbf{P} to obtain $\hat{\mathbf{I}}''_s$. Finally, we add a reconstruction loss: $\mathcal{L}_{\text{CVC}} = \lambda_{\text{CVC}} \mathcal{L}_s(\hat{\mathbf{I}}''_s, \hat{\mathbf{I}}_s)$.

the pixel-aligned conditional discriminator D_s concatenates color images and label maps as input, which encourages pixel alignment between color images and label maps. Notably, in D_s , we stop the gradients for the color images to prevent a potential quality downgrade. We also feed the rendered low-res images to prevent the upsampler from hallucinating details, inconsistent with the low-res output. The adversarial loss can be written as follows.

$$\mathcal{L}_{\text{GAN}} = \lambda_{D_c} \mathcal{L}_{D_c}(\hat{\mathbf{I}}_c^+, \hat{\mathbf{I}}_c) + \lambda_{D_s} \mathcal{L}_{D_s}(\hat{\mathbf{I}}_c^+, \hat{\mathbf{I}}_c, \hat{\mathbf{I}}_s^+, \hat{\mathbf{I}}_s).$$

where λ_{D_c} and λ_{D_s} balance two terms. To stabilize the GAN training, we adopt the R1 regularization loss [45].

Cross-view Consistency Loss. We observe that inputting label maps of the same object from different viewpoints will sometimes result in different 3D shapes. Therefore we add a cross-view consistency loss to regularize the training, as illustrated in Figure 3. Given an input label map \mathbf{I}_s and its associated pose \mathbf{P} , we generate the label map $\hat{\mathbf{I}}_s$ from a different viewpoint \mathbf{P}' , and render the label map $\hat{\mathbf{I}}_s$ back to the pose \mathbf{P} using $\hat{\mathbf{I}}_s$ as input. We add a reconstruction loss between $\hat{\mathbf{I}}_s$ and $\hat{\mathbf{I}}_s$:

$$\mathcal{L}_{\text{CVC}} = \lambda_{\text{CVC}} \mathcal{L}_s(\hat{\mathbf{I}}_s'', \hat{\mathbf{I}}_s),$$

where \mathcal{L}_s denotes the reconstruction loss in the label space, and λ_{CVC} weights the loss term. This loss is crucial for reducing error accumulation during cross-view editing.

Optimization. Our final learning objective is written as follows:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{recon}} + \mathcal{L}_{\text{GAN}} + \mathcal{L}_{\text{CVC}}.$$

At every iteration, we determine whether to use a ground-truth pose or sample a random one with a probability of p . We use the reconstruction loss and GAN loss for ground-truth poses, while for random poses, we only use the GAN loss. We provide the hyper-parameters and more implementation details in Appendix B.



Figure 4. **Qualitative Comparison with Pix2NeRF [6], SoFGAN [11], and SEAN [89]** on CelebAMask dataset for seg2face task. SEAN fails in multi-view synthesis, while SoFGAN suffers from multi-view inconsistency (e.g., face identity changes across viewpoints). Our method renders high-quality images while maintaining multi-view consistency. Please check our [website](#) for more examples.

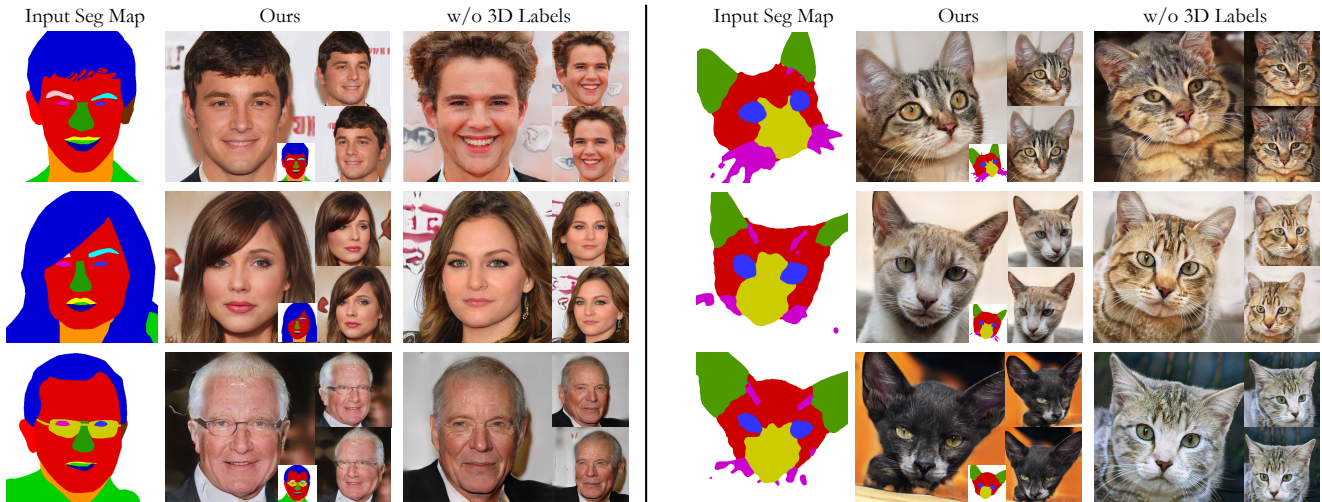


Figure 5. **Qualitative ablation on seg2face and seg2cat.** We ablate our method by removing the branch that renders label maps (*w/o 3D Labels*). Our results better align with input labels (e.g., hairlines and the cat’s ear).

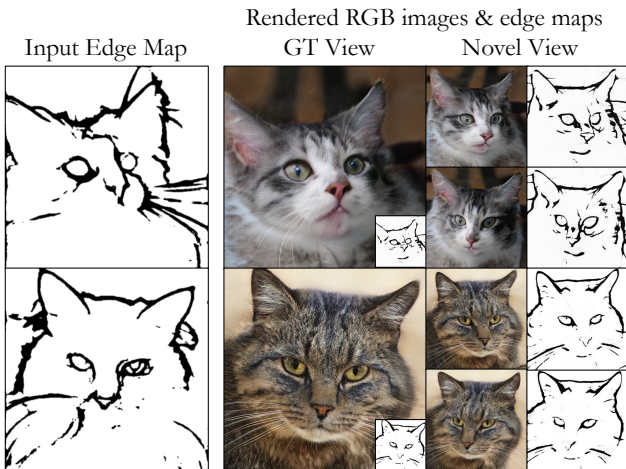


Figure 6. **Results on edge2cat.** Our model is trained on AFHQ-cat [16] with edges extracted by pidinet [66].

4. Experiment

We first introduce the datasets and evaluation metrics. Then we compare our method with the baselines. Finally, we demonstrate cross-view editing and multi-modal synthesis applications enabled by our method.

Datasets. We consider four tasks: *seg2face*, *seg2cat*,

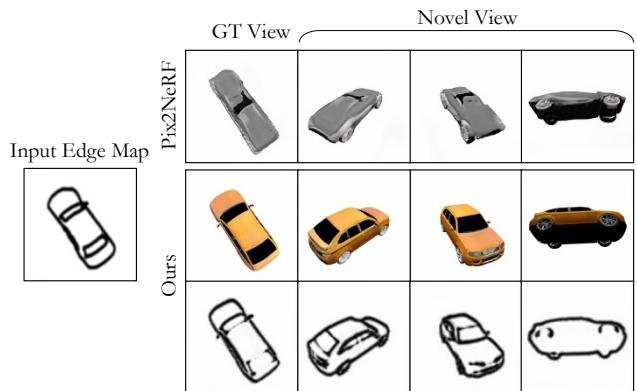


Figure 7. **Qualitative comparisons on edge2car.** *pix2pix3D* (Ours) and Pix2NeRF [6] are trained on shapenet-car [10], and *pix2pix3D* achieves better quality and alignment than Pix2NeRF.

edge2cat, and *edge2car* in our experiments. For *seg2face*, we use CelebAMask-HQ [38] for evaluation. CelebAMask-HQ contains 30,000 high-resolution face images from CelebA [41], and each image has a facial part segmentation mask and a predicted pose. The segmentation masks contain 19 classes, including skin, eyebrows, ears, mouth, lip, etc. The pose associated with each image segmentation is predicted by HopeNet [60]. We split the CelebAMask-HQ dataset into

Seg2Face	QUALITY			ALIGNMENT		FVV Identity ↓
	FID ↓	KID ↓	SG Diversity ↑	mIoU ↑	acc ↑	
CELEBAMASK [38]						
SEAN [89]	32.74	0.018	0.29	0.52	0.85	N/A
SoFGAN [11]	23.34	0.012	0.33	0.53	0.89	0.58
PIX2NeRF [6]	54.23	0.042	0.16	0.36	0.65	0.44
PIX2PIX3D (OURS)						
w/o 3D LABELS	12.96	0.005	0.30	N/A (0.43)	N/A (0.81)	0.38
w/o CVC	11.62	0.004	0.30	0.50 (0.50)	0.87 (0.85)	0.42
FULL MODEL	11.54	0.003	0.28	0.51 (0.52)	0.90 (0.88)	0.36
FULL MODEL†	11.13	0.003	0.29	0.51 (0.50)	0.90 (0.87)	0.36

Table 1. **Seg2face Evaluation.** Our metrics include image quality (FID, KID, SG Diversity), alignment (mIoU and acc against GT label maps), and multi-view consistency (FVV Identity). Single-generation diversity (SG Diversity) is obtained by computing the LPIPS metric between randomly generated pairs given a single conditional input. To evaluate alignment, we compare the generated label maps against the ground truth in terms of mIoU and pixel accuracy (acc). Alternatively, given a generated image, one could estimate label maps via a face parser, and compare those against the ground truth (numbers in parentheses). We include SEAN [89] and SoFGAN [11] as baselines, and modify Pix2NeRF [6] to take conditional input. Our method achieves the best quality, alignment ACC, and FVV Identity while being competitive on SG Diversity. SoFGAN tends to have better alignment but worse 3D consistency. We also ablate our method w.r.t the 3D labels and the cross-view consistency (CVC) loss. Our 3D labels are crucial for alignment, while the CVC loss improves multi-view consistency. Using pre-trained models from EG3D (†) also improves the performance.

Edge2Car	QUALITY			ALIGNMENT
	FID ↓	KID ↓	SG Diversity ↑	AP ↑
PIX2NeRF [6]	23.42	0.014	0.06	0.28
PIX2PIX3D (OURS)				
w/o 3D LABELS	10.73	0.005	0.12	0.45 (0.42)
w/o CVC	9.42	0.004	0.13	0.61 (0.59)
FULL MODEL	8.31	0.004	0.13	0.63 (0.59)

Table 2. **Edge2car Evaluation.** We compare our method with Pix2NeRF [6] on edge2car using the shapenet-car [10] dataset. Similar to Table 1, we evaluate FID, KID, and SG Diversity for image quality. We also evaluate the alignment with the input edge map using AP. Similarly, we can either run informative drawing [7] on generated images to obtain edge maps (numbers in parentheses) or directly use generated edge maps to calculate the metrics. We achieve better image quality and alignment than Pix2NeRF. We also find that using 3D labels and cross-view consistency loss is helpful regarding FID and AP metrics.

a training set of 24,183, a validation set of 2,993, and a test set of 2,824, following the original work [38]. For seg2cat and edge2cat, we use AFHQ-cat [16], which contains 5,065 images at 512× resolution. We estimate the viewpoints using unup3d [75]. We extract the edges using pidinet [66] and obtain segmentation by clustering DINO features [2] into 6 classes. For edge2car, we use 3D models from shapenet-

Seg2Cat	QUALITY			ALIGNMENT	
	FID ↓	KID ↓	SG Diversity ↑	mIoU ↑	acc ↑
AFHQ-CAT [34]					
PIX2NeRF [6]	43.92	0.081	0.15	0.27	0.58
OURS					
w/o 3D LABELS	10.41	0.004	0.26	N/A (0.49)	N/A (0.69)
w/o CVC	9.64	0.004	0.26	0.66 (0.63)	0.76 (0.73)
FULL MODEL	8.62	0.003	0.27	0.66 (0.62)	0.78 (0.73)

Table 3. **Seg2cat Evaluation.** We compare our method with Pix2NeRF [6] on Seg2Cat using AFHQ-cat dataset [16], with segmentation obtained by clustering DINO features [2]. Similar to Table 1, we evaluate the image quality and alignment. Ours performs better in all metrics.

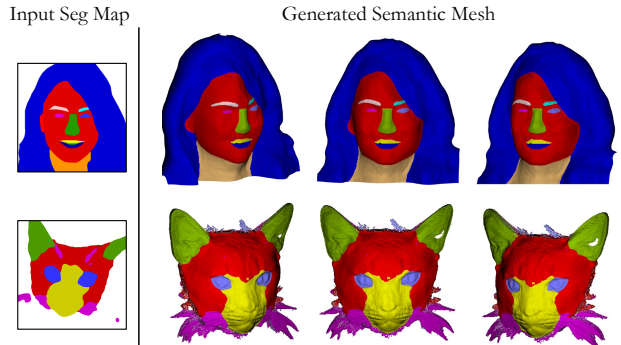


Figure 8. **Semantic Mesh.** We show semantic meshes of human and cat faces from marching cubes colored by 3D labels.

car [10] and render 500,000 images at 128× resolution for training, and 30,000 for evaluation. We extract the edges using informative drawing [7]. We train our model at 512× resolution except for 128× in the edge2car task.

Running Time. For training the model at 512× resolution, it takes about three days on eight RTX 3090 GPUs. But we can significantly reduce the training time to 4 hours if we initialize parts of our model with pretrained weights from EG3D [8]. During inference, our model takes 10 ms to obtain the style vector, and another 30 ms to render the final image and the label map on a single RTX A5000. The low latency (25 FPS) allows for interactive user editing.

4.1. Evaluation metrics

We evaluate the models from two aspects: 1) the image quality regarding fidelity and diversity, and 2) the alignment between input label maps and generated outputs.

Quality Metrics. Following prior works [21, 57], we use the clean-fid library [58] to compute Fréchet Inception Distance (FID) [23] and Kernel Inception Distance (KID) [4] to measure the distribution distance between synthesized results and real images. We also evaluate the single-generation diversity (SG Diversity) by calculating the LPIPS metric between randomly generated pairs given a single input following prior works [11, 87]. For FID and KID, we generate

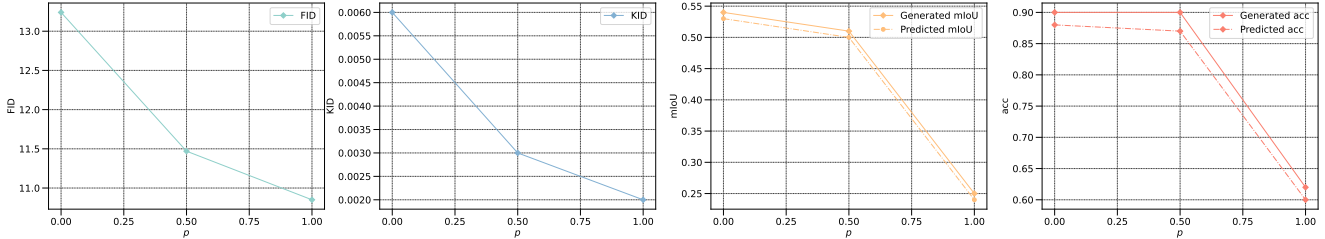


Figure 9. We study the effect of random pose sampling probability p during training. Without random poses ($p = 0$), the model achieves the best alignment with input semantic maps, with reduced image quality. In contrast, *only* using random poses ($p = 1$) achieves the best image quality, while results fail to align with input maps. We find $p = 0.5$ balances the image quality and input alignment.

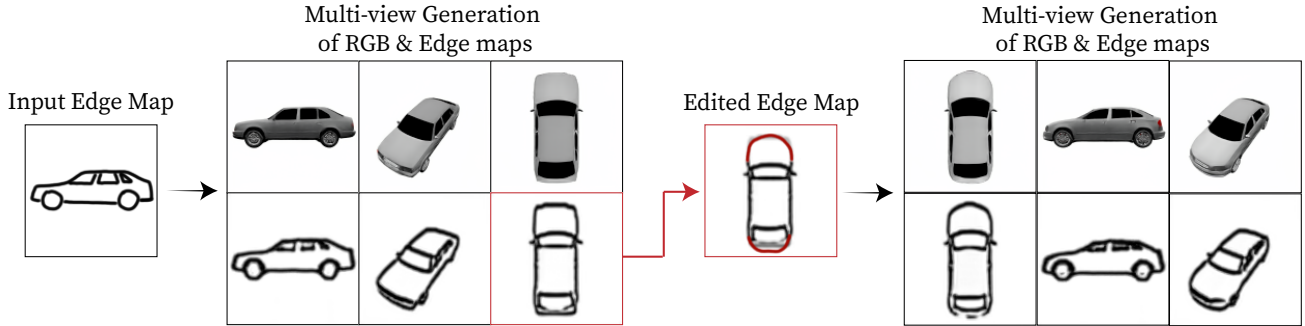


Figure 10. **Cross-view Editing of Edge2Car.** Our 3D editing system allows users to edit label maps from any viewpoint instead of only the input view. Importantly, our feed-forward encoder allows fast inference of the latent code without GAN-inversion. Typically, a single forward pass of rendering takes only 40 ms on a single RTX A5000, which enables interactive editing. Please check our demo video on our [website](#).

10 images per label map in the test set using randomly sampled z . We compare our generated images with the whole dataset, including training and test images.

Alignment Metrics. We evaluate models on the test set using mean Intersection-over-Union (mIoU) and pixel accuracy (acc) for segmentation maps following existing works [57, 63], and average precision (AP) for edge maps. For those models that render label maps as output, we directly compare them with ground-truth labels. Otherwise, we first predict the label maps from the output RGB images using off-the-shelf networks [38, 66], and then compare the prediction with the ground truth. The metrics regarding such predicted semantic maps are reported within brackets in Table 1 and Table 2.

For seg2face, we evaluate the preservation of facial identity from different viewpoints (FVV Identity) by calculating their distances with the dlib face recognition algorithm*.

4.2. Baseline comparison

Baselines. Since there are no prior works on conditional 3D-aware image synthesis, we make minimum modifications to Pix2NeRF [6] to be conditional on label maps instead of images. For a thorough comparison, we introduce several baselines: SEAN [89] and SoFGAN [11]. 2D baselines like SEAN [89] cannot generate multi-view images by design (N/A for FVV Identity), while SoFGAN [11] uses

an unconditional 3D semantic map generator before the 2D generator so we can evaluate FVV Identity for that.

Results. Figure 4 shows the qualitative comparison for seg2face and Table 1 reports the evaluation results. SoFGAN [11] tends to produce results with slightly better alignment but worse 3D consistency for its 2D RGB generator. Our method achieves the best quality, alignment acc, and FVV Identity while being competitive with 2D baselines on SG diversity. Figure 5 shows the qualitative ablation on seg2face and seg2cat. Table 5 reports the metrics for seg2cat. Figure 6 shows the example results for edge2cat. Figure 7 shows the qualitative comparison for edge2car and Table 2 reports the metrics. Our method achieves the best image quality and alignment. Figure 8 shows semantic meshes of human and cat faces, extracted by marching cubes and colored by our learned 3D labels. We provide more evaluation results in Appendix A.

Ablation Study. We compare our full method to several variants. Specifically, (1) w/o 3D LABELS, we remove the branch of rendering label maps from our method, and (2) w/o CVC, we remove the cross-view consistency loss. From Table 1, Table 2, and Figure 5, rendering label maps is crucial for the alignment with the input. We posit that the joint learning of appearance, geometry, and label information poses strong constraints on correspondence between the input label maps and the 3D representation. Thus our method can synthesize images pixel-aligned with the inputs. Our CVC loss helps preserve the facial identity from different viewpoints.

*https://github.com/ageitgey/face_recognition

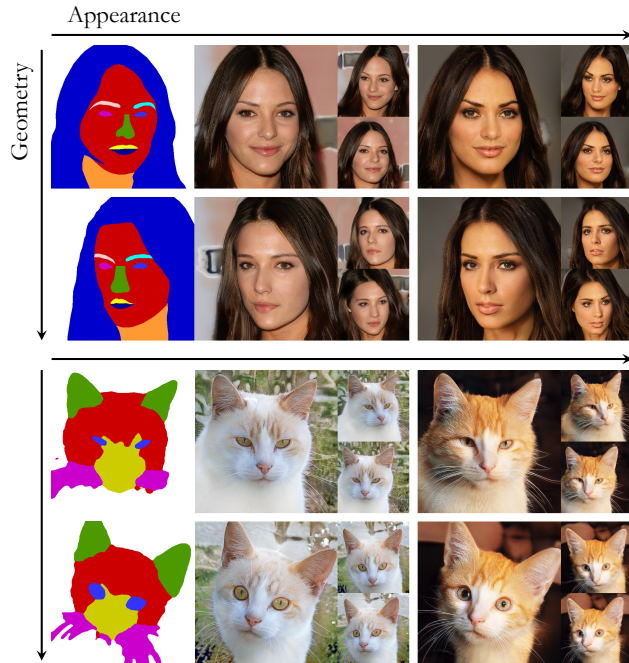


Figure 11. **Multi-modal Synthesis.** The leftmost column is the input segmentation map. We use the same segmentation map for each row. We generate multi-modal results by randomly sampling an appearance style for each column.

Analysis on random sampling of poses. We study the effect of the different probabilities of sampling random poses during training, as shown in Figure 9. When sampling no random poses ($p = 0$), the model best aligns with input label maps with suboptimal image quality. Conversely, *only* sampling random poses ($p = 1$) gives the best image quality but suffers huge misalignment with input label maps. We find $p = 0.5$ achieves the balance between the image quality and the alignment with the input.

4.3. Applications

Cross-view Editing. As shown in Figure 10, our 3D editing system allows users to generate and edit label maps from any viewpoint instead of only the input view. The edited label map is further fed into the conditional encoder to update the 3D representation. Unlike GAN inversion [85], our feed-forward conditional encoder allows fast inference of the latent code. Thus, a single forward pass of our full model takes only 40 ms on a single RTX A5000.

Multi-modal synthesis and interpolation. Like other style-based generative models [8, 21, 34, 36], our method can disentangle the geometry and appearance information. Specifically, the input label map captures the geometry information while the randomly sampled latent code controls the appearance. We show style manipulation results in Figure 11. We can also interpolate both the geometry styles and the appearance styles (Figure 12). These results show the clear disentanglement of our 3D representation.

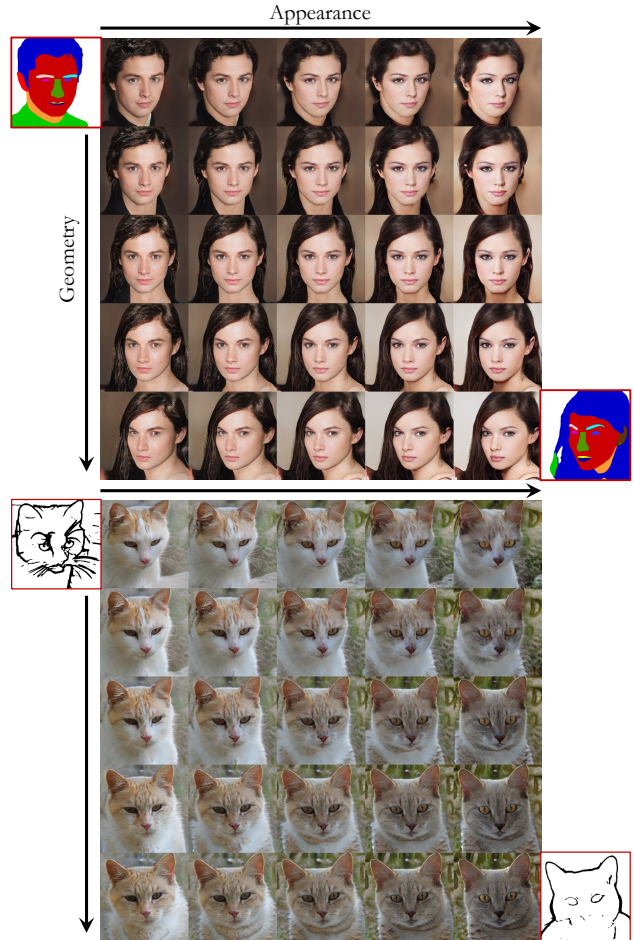


Figure 12. **Interpolation.** In each 5×5 grid, the images at the top left and bottom right are generated from the input maps next to them. Each row interpolates two images in label space, while each column interpolates the appearance. For camera poses, we interpolate the pitch along the row and the yaw along the column.

5. Discussion

We have introduced `pix2pix3D`, a 3D-aware conditional generative model for controllable image synthesis. Given a 2D label map, our model allows users to render images given any viewpoint. Our model augments the neural field with 3D labels, assigning label, color, and density to every 3D point, allowing for the simultaneous rendering of the image and a pixel-aligned label map. The learned 3D labels further enable interactive 3D cross-view editing. We discuss the broader impact and limitations in the appendix.

Acknowledgments. We thank Sheng-Yu Wang, Nupur Kumari, Gaurav Parmar, Ruihan Gao, Muyang Li, George Cazenavette, Andrew Song, Zhipeng Bao, Tamaki Kojima, Krishna Wadhvani, Takuya Narihira, and Tatsuo Fujiwara for their discussion and help. We are grateful for the support from Sony Corporation, Singapore DSTA, and the CMU Argo AI Center for Autonomous Vehicle Research.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [2] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. *ECCVW What is Motion For?*, 2022. 6
- [3] David Bau, Hendrik Strobelt, William Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. In *ACM SIGGRAPH*, 2019. 2
- [4] Mikolaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations (ICLR)*, 2018. 6
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations (ICLR)*, 2019. 2
- [6] Shengqu Cai, Anton Obukhov, Dengxin Dai, and Luc Van Gool. Pix2nerf: Unsupervised conditional π -gan for single image to neural radiance fields translation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 5, 6, 7, 13
- [7] Caroline Chan, Frédo Durand, and Phillip Isola. Learning to generate line drawings that convey geometry and semantics. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 6
- [8] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 4, 6, 8, 14
- [9] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [10] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. 2, 5, 6
- [11] Anpei Chen, Ruiyang Liu, Ling Xie, Zhang Chen, Hao Su, and Jingyi Yu. Sofgan: A portrait image generator with dynamic styling. In *ACM SIGGRAPH*, 2021. 2, 5, 6, 7, 13, 14
- [12] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [13] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [14] Tao Chen, Zhe Zhu, Ariel Shamir, Shi-Min Hu, and Daniel Cohen-Or. 3-sweep: Extracting editable objects from a single photo. *ACM Transactions on Graphics (TOG)*, 32(6):1–10, 2013. 2
- [15] Yuedong Chen, Qianyi Wu, Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Sem2nerf: Converting single-view semantic masks to neural radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [16] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5, 6, 14
- [17] Johanna Delanoy, Adrien Bousseau, Mathieu Aubry, Phillip Isola, and Alexei A Efros. What you sketch is what you get: 3d sketching using multi-view deep volumetric prediction. In *ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games (I3D)*, 2018. 2
- [18] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [19] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1
- [20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014. 1, 2, 4, 13
- [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 (ICLR)*, 2022. 2, 3, 4, 6, 8, 13, 14
- [22] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato’s cave: 3d shape from adversarial rendering. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [23] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 6
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1
- [25] Hsin-Ping Huang, Hung-Yu Tseng, Hsin-Ying Lee, and Jia-Bin Huang. Semantic view synthesis. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [26] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *European Conference on Computer Vision (ECCV)*, 2018. 2
- [27] Zeng Huang, Tianye Li, Weikai Chen, Yajie Zhao, Jun Xing, Chloe Legendre, Linjie Luo, Chongyang Ma, and Hao Li. Deep volumetric video from very sparse multi-view performance capture. In *European Conference on Computer Vision (ECCV)*, pages 351–369, 2018. 2

- [28] Takeo Igarashi, Satoshi Matsuoka, and Hidehiko Tanaka. Teddy: a sketching interface for 3d freeform design. In *ACM SIGGRAPH*, 1999. 2
- [29] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2
- [30] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [31] Kaiwen Jiang, Shu-Yu Chen, Feng-Lin Liu, Hongbo Fu, and Lin Gao. Nerffacediting: Disentangled face editing in neural radiance fields. In *ACM SIGGRAPH Asia*, 2022. 2
- [32] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH*, 18(3):165–174, 1984. 3
- [33] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2
- [34] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2, 6, 8
- [35] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [36] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 8, 13, 14
- [37] Natasha Kholgade, Tomas Simon, Alexei Efros, and Yaser Sheikh. 3d object manipulation in a single photograph using stock 3d models. *ACM Transactions on Graphics (TOG)*, 33(4):1–12, 2014.
- [38] 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. 2, 5, 6, 7, 14
- [39] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [40] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. *Advances in neural information processing systems*, 30, 2017. 2
- [41] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. 5
- [42] Zhaoliang Lun, Matheus Gadelha, Evangelos Kalogerakis, Subhransu Maji, and Rui Wang. 3d shape reconstruction from sketches via multi-view convolutional networks. In *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017. 2
- [43] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [44] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [45] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International Conference on Machine Learning (ICML)*, 2018. 4, 13
- [46] 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 (ECCV)*, 2020. 2, 3
- [47] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2
- [48] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. In *ACM SIGGRAPH*, 2022. 2
- [49] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [50] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [51] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [52] 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. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 4
- [53] Xingang Pan, Xudong Xu, Chen Change Loy, Christian Theobalt, and Bo Dai. A shading-guided generative implicit model for shape-accurate 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2
- [54] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [55] Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. In *European Conference on Computer Vision (ECCV)*, 2020. 2

- [56] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [57] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 6, 7
- [58] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 6
- [59] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [60] Nataniel Ruiz, Eunji Chong, and James M. Rehg. Fine-grained head pose estimation without keypoints. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshop*, 2018. 5
- [61] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [62] Edgar Schönfeld, Vadim Sushko, Dan Zhang, Juergen Gall, Bernt Schiele, and Anna Khoreva. You only need adversarial supervision for semantic image synthesis. In *International Conference on Learning Representations (ICLR)*, 2020. 2
- [63] Edgar Schönfeld, Vadim Sushko, Dan Zhang, Juergen Gall, Bernt Schiele, and Anna Khoreva. You only need adversarial supervision for semantic image synthesis. In *International Conference on Learning Representations (ICLR)*, 2021. 7
- [64] 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
- [65] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [66] Zhuo Su, Wenzhe Liu, Zitong Yu, Dewen Hu, Qing Liao, Qi Tian, Matti Pietikainen, and Li Liu. Pixel difference networks for efficient edge detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 5, 6, 7
- [67] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew Davison. iMAP: Implicit mapping and positioning in real-time. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 2
- [68] 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. In *ACM Transactions on Graphics (TOG)*, 2022. 2
- [69] Jingxiang Sun, Xuan Wang, Yong Zhang, Xiaoyu Li, Qi Zhang, Yebin Liu, and Jue Wang. Fenerf: Face editing in neural radiance fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 4
- [70] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [71] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [72] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. In *ACM Transactions on Graphics (TOG)*, 2021. 2
- [73] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot...for now. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 14
- [74] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [75] Shangzhe Wu, Christian Ruppert, and Andrea Vedaldi. Un-supervised learning of probably symmetric deformable 3d objects from images in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6
- [76] Xiaohua Xie, Kai Xu, Niloy J Mitra, Daniel Cohen-Or, Wenyong Gong, Qi Su, and Baoquan Chen. Sketch-to-design: Context-based part assembly. In *Computer Graphics Forum*, volume 32, pages 233–245. Wiley Online Library, 2013. 2
- [77] Yinghao Xu, Sida Peng, Ceyuan Yang, Yujun Shen, and Bolei Zhou. 3d-aware image synthesis via learning structural and textural representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [78] Shunyu Yao, Tzu Ming Hsu, Jun-Yan Zhu, Jiajun Wu, Antonio Torralba, Bill Freeman, and Josh Tenenbaum. 3d-aware scene manipulation via inverse graphics. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2
- [79] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. Pixelnerf: Neural radiance fields from one or few images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [80] Jichao Zhang, Enver Sangineto, Hao Tang, Aliaksandr Siarohin, Zhun Zhong, Nicu Sebe, and Wei Wang. 3d-aware semantic-guided generative model for human synthesis. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [81] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [82] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 4
- [83] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 13

- [84] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image editing. In *European Conference on Computer Vision (ECCV)*, 2020. [2](#)
- [85] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *European Conference on Computer Vision (ECCV)*, 2016. [2](#), [8](#)
- [86] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. [1](#), [2](#)
- [87] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. *Advances in neural information processing systems*, 30, 2017. [6](#)
- [88] Jun-Yan Zhu, Zhoutong Zhang, Chengkai Zhang, Jiajun Wu, Antonio Torralba, Josh Tenenbaum, and Bill Freeman. Visual object networks: Image generation with disentangled 3d representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. [2](#)
- [89] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. Sean: Image synthesis with semantic region-adaptive normalization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#), [5](#), [6](#), [7](#)
- [90] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R. Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [2](#)

Appendix

We include additional experimental results, implementation details, and the societal impact of our work. Please also view our [webpage](#) for our interactive editing demo video and additional visual results.

A. Additional Experiments

Cross-view Editing of Seg2cat. In addition to the edge2car editing example in the main paper, we showcase the editability of segmentation maps in Figure 13. Note that the edited segmentation map does not have to be from the same viewpoint as the input segmentation map.

Ablation study on discriminator design. In the main paper, we introduce our Pixel-aligned Conditional Discriminator that concatenates RGB images and label maps as input. To verify the effectiveness of our discriminator design, we introduce three ablation experiments in Table 4. We find that our image discriminator helps improve the image quality, while our pixel-aligned conditional discriminator is crucial for the alignment.

CELEBA-MASK	QUALITY		ALIGNMENT	
	FID ↓	KID ↓	mIoU ↑	acc ↑
OURS	11.54	0.003	0.51 (0.52)	0.90 (0.88)
W/O IMAGE D	15.32	0.006	0.51 (0.52)	0.89 (0.85)
W/O CONDITIONAL D	12.02	0.004	0.37 (0.47)	0.82 (0.80)
W/O PIXEL-ALIGN D	11.94	0.003	0.41 (0.40)	0.82 (0.81)

Table 4. **Ablation Study on Discriminator Design.** To verify the effectiveness of our discriminator design, we introduce three ablation experiments: (1)W/O IMAGE D, we remove the image discriminator and only keep the conditional discriminator that accepts the concatenation of image and segmentation maps; (2)W/O CONDITIONAL D, we remove the conditional discriminator and only keeps the image discriminator; (3)W/O PIXEL-ALIGN D, we keep both discriminators, but the conditional discriminator no longer concatenates color images as part of the input. Our image discriminator improves the image quality, while our pixel-aligned conditional discriminator ensures alignment.

Evaluation on Seg2Car. We evaluate our method on an additional non-face dataset Seg2Car, where we get the segmentation model from DatasetGAN [83]. We show the visual and evaluation results in the Figure 14 and Table 5. We find our method outperforms Pix2NeRF [6].

Figure 15 compares our method with SoFGAN [11] regarding multi-view consistency. We also show our method’s capability of correcting errors in the user input in Figure 16.

Seg2Car	QUALITY			ALIGNMENT	
	SHAPNET-CAR FID ↓	KID ↓	SG Diversity ↑	mIoU ↑	acc ↑
PIX2NERF	25.86	0.018	0.08	0.24	0.59
OURS	9.35	0.004	0.14	0.58	0.88

Table 5. **Seg2car Evaluation.** We compare our method with Pix2NeRF [6]. Ours performs better in all metrics.

B. Implementation Details

Class-balanced cross-entropy. In Section 3.2 of the main text, we mentioned using class-balanced cross-entropy loss for reconstructing 2D segmentation maps. Specifically,

$$\mathcal{L}_s(\hat{\mathbf{I}}_s, \mathbf{I}_s^+) = \mathbb{E}_n - \sum_{c=1}^C w_c \log \frac{\exp(x_{n,c})}{\sum_{i=1}^C \exp(x_{n,i})} y_{n,c}.$$

where $x_{n,c}$ is the semantic logits of class c at location n of \mathbf{I}_s^+ , $y_{n,c}$ is the ground-truth probability of class c at location n of $\hat{\mathbf{I}}_s$, and w_c is the weight of each class c .

In our case, 2D segmentation maps are imbalanced as skin and hair cover a lot more areas than the other classes. So we calculate w_c based on the inverse frequency of the classes in the training set,

$$w_c = \sqrt{\frac{\# \text{ of pixels with class } c}{\# \text{ of all pixels}}}.$$

Regularization. As mentioned in the main text, we use non-saturating loss [20] and R1 Regularization [45] for GAN training following [21, 36]. Specifically,

$$\mathcal{L}_{\text{GAN}}(G, D) = \mathbb{E}[f(D(G(z), \hat{\mathbf{I}}_s))] + \mathbb{E}[f(-D(\hat{\mathbf{I}}_c)) + \lambda \|\nabla D(\hat{\mathbf{I}}_c)\|^2],$$

where G is the generator, D is the discriminator, $f(u) = -\log(1 + \exp(-u))$, and $\lambda = 0.5$.

Hyper-parameters. $\lambda_c = 1$, $\lambda_s = 5$ for edge maps, $\lambda_s = 1$ for segmentation maps, $\lambda_{D_c} = 1$, $\lambda_{D_s} = 0.1$, and $\lambda_{\text{CVC}} = 1e - 5$. Check our [codebase](#) for more detailed hyperparameters.

C. Discussion

Broader Impact. Our work allows a novice user to create 3D content more easily. The 3D outputs can be directly used in photo editing software as well as virtual reality and augmented reality applications.

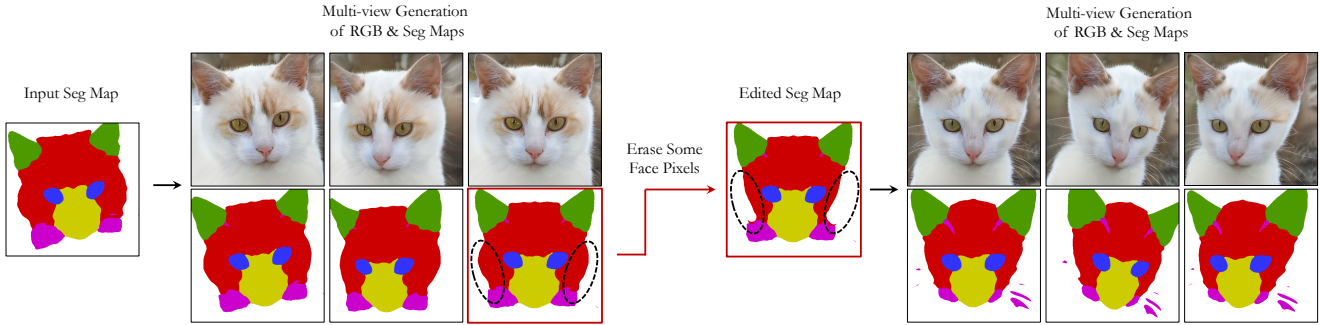


Figure 13. **Cross-view Editing of Seg2cat.** The 3D representation can be edited from a viewpoint different than the input seg map.

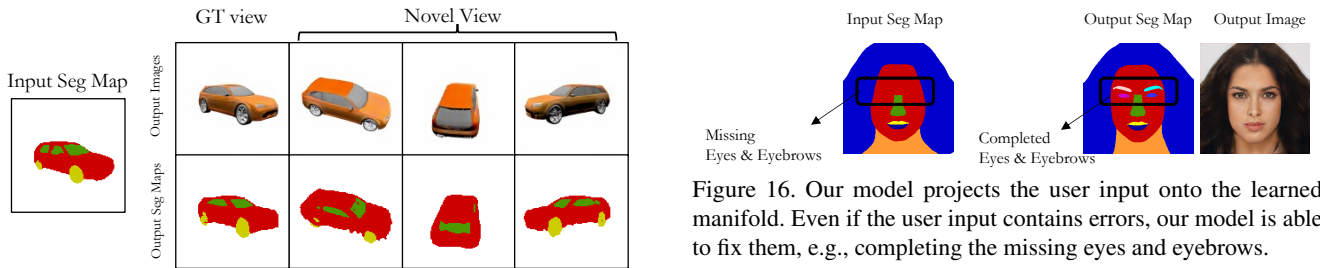


Figure 14. **Visual Results of Seg2Car.**

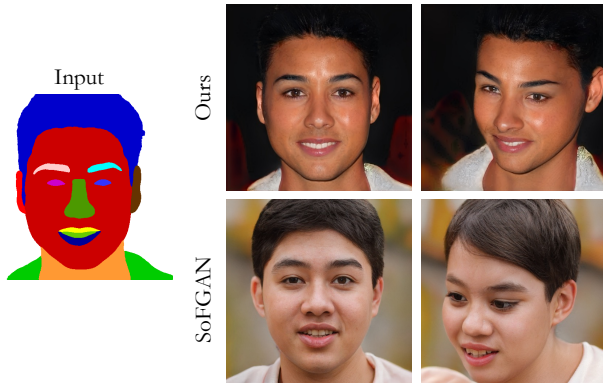


Figure 15. **Multi-view Consistency.** We compare our method with SoFGAN [11] regarding multi-view consistency. Although SoFGAN can generate images from different viewpoints, their method shifts the face identity across viewpoints. In contrast, our method better preserves the identity.

Similar to recent works on data-driven 2D and 3D face synthesis, we suffer from biases in the underlying dataset. Our model is trained on CelebAMask-HQ dataset, as it provides segmentation masks that can be used as conditional input. To reduce the dataset bias, one future direction is to run our model on more diverse datasets with a pre-trained face parser. While our work allows for controllable 3D content generation, there may be potential misuse of the generated content. As an attempt to identify the generated content from the real photos, we run a forensics detector [73] on our gen-

erated results, and find our generated images can be detected with an accuracy of 89.77%, and an average precision of 99.97%.

Usage of Existing Assets. We use CelebAMask-HQ dataset [38]. The CelebA dataset is available for non-commercial research purposes only. All images of the CelebA dataset are obtained from the Internet. The face identities are released upon request for research purposes only. See [CelebA website](#) for details. We also use AFHQ Cat dataset [16]. This dataset is under [Creative Commons license](#). Our work is also inspired by a few codebases. StyleNeRF codebase [21] is under Creative Common license. StyleGAN2 codebase [36] is under the [Nvidia Source Code License](#). EG3D codebase [8] is under the [Nvidia Source Code License](#).

Limitations. Our current method has three major limitations. First, it mainly focuses on modeling the appearance and geometry of a single object category. Extending the method to more complex scene datasets with multiple objects is a promising next step, though defining a canonical pose for generic scenes poses a nontrivial challenge. Second, our model’s generation is limited to the dataset’s distribution. Our model will not follow the user input unless it is within the dataset’s distribution. Incorporating diffusion models and training on more diverse datasets can potentially improve the generalization. Finally, our model training requires camera poses associated with each training image, though our method does not require poses during inference time. Eliminating the requirement for pose information will further broaden the scope of applications.

D. Changelog

V2. Add more citations in Section 2. Fix some typos.

V1. Initial preprint release (CVPR 2023).