


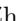




Sem2NeRF: Converting Single-View Semantic Masks to Neural Radiance Fields

Yuedong Chen¹, Qianyi Wu¹, Chuanxia Zheng¹,
Tat-Jen Cham², and Jianfei Cai¹

¹ Monash University, Australia

{yuedong.chen, qianyi.wu, jianfei.cai}@monash.edu

² Nanyang Technological University, Singapore

chuanxia001@e.ntu.edu.sg, astjcham@ntu.edu.sg

Abstract. Image translation and manipulation have gain increasing attention along with the rapid development of deep generative models. Although existing approaches have brought impressive results, they mainly operated in 2D space. In light of recent advances in NeRF-based 3D-aware generative models, we introduce a new task, Semantic-to-NeRF translation, that aims to reconstruct a 3D scene modelled by NeRF, conditioned on one single-view semantic mask as input. To kick-off this novel task, we propose the Sem2NeRF framework. In particular, Sem2NeRF addresses the highly challenging task by encoding the semantic mask into the latent code that controls the 3D scene representation of a pre-trained decoder. To further improve the accuracy of the mapping, we integrate a new region-aware learning strategy into the design of both the encoder and the decoder. We verify the efficacy of the proposed Sem2NeRF and demonstrate that it outperforms several strong baselines on two benchmark datasets. Code and video are available at <https://donydchen.github.io/sem2nerf/>.

Keywords: NeRF-based generation, conditional generative model, 3D deep learning, neural radiance fields, image-to-image translation

1 Introduction

Controllable image generation, translation, and manipulation have seen rapid advances in the last few years along with the emergence of Generative Adversarial Networks (GANs) [14]. Current systems are able to freely change the image appearance through referenced images [21,70,20], modify scene content via semantic masks [53,40,29], and even accurately manipulate various attributes in feature space [24,56,57]. Despite impressive performance and wide applicability, these systems are mainly focused on 2D images, without directly considering the 3D nature of the world and the objects within.

Concurrently, significant progress has been made for 3D generation by using deep generative networks [14,26]. Methods were developed for different 3D shape representations, including voxels [54], point clouds [36], and meshes [13].

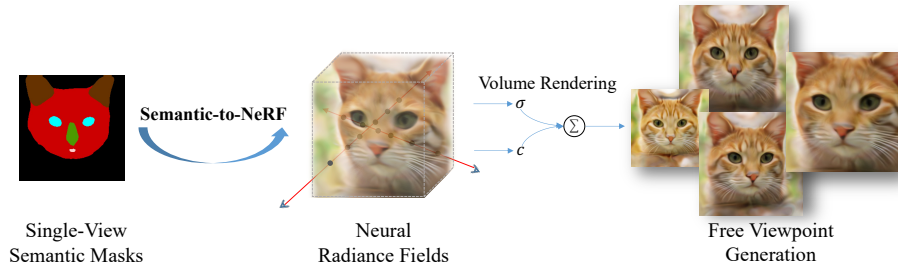


Fig. 1. Illustration of the Semantic-to-NeRF translation task, which aims to achieve free-viewpoint image generation by taking only a single-view semantic mask as input

More recently, Neural Radiance Fields (NeRF) [38] has been a new paradigm for 3D representation, providing accurate 3D shape and view-dependent appearance simultaneously. Based on this new representation, seminal 3D generation approaches [46,39,3,15,2] have been proposed that aim to generate photorealistic images from a given distribution in a 3D-aware and view-consistent manner. However, these techniques are primarily developed purely for high-quality 3D generation, leaving controllable 3D manipulation and editing unsolved.

It would be a dramatic enhancement if we can *freely manipulate and edit an object’s content and appearance in 3D space, while only leveraging easily obtained 2D input information*. In this paper, we take an initial step toward this grand goal by introducing a new task, termed **Semantic-to-NeRF translation**, analogous to a 2D Semantic-to-Image translation task but operating on 3D space. Specifically, Semantic-to-NeRF translation (see Fig. 1) takes as input a single-view 2D semantic mask, yet output a NeRF-based 3D representation that can be used to render photorealistic images in a 3D-aware view-consistent manner. More importantly, it allows free editing of the object’s content and appearance in 3D space, by modifying the content only via a single-view 2D semantic mask.

However, generating 3D structure from a single 2D image is already an ill-posed problem, and it will be even more so from a single 2D semantic mask. There are also two other major issues in this novel task:

1. *Large information gap between 3D structure and 2D semantics.* A single-view 2D semantic mask neither holds any 3D shape or surface information, nor provides much guidance for plausible appearances, making it tough to generate a neural radiance field with comprehensive details.
2. *Imbalanced semantic distribution.* Since semantic classes tend to be area-imbalanced within an image, *e.g.* eyes occupy less than 1% of a face while hair can take up larger than 40%, existing CNN-based networks may over-attend to larger semantic regions, while discounting smaller semantic regions that may be perceptually more salient. This will result in poor controllable editing in 3D space when we alter small semantic regions.

To mitigate these issues, we propose a novel framework, **Sem2NeRF**, that builds on NeRF [38] for 3D representation, by augmenting it with a seman-

tic translation branch that conditionally generates high-quality 3D-consistent images. In particular, the framework is based on an encoder-decoder architecture that converts a single-view 2D semantic mask to an embedded code, and then transfers it to a NeRF representation for rendering 3D-consistent images.

Our broad idea here is that, instead of directly learning to predict 3D structure from degenerate single-view 2D semantic masks, *the network can alternatively learn the 3D shape and appearance representation from large numbers of unstructured 2D RGB images*. This has achieved significant advances in NeRF-based generator [46,39,3,15,2], which transforms a random vector to a NeRF representation. In short, our scenario is thus: we have a well-trained 3D generator, but we aim to further control the generated content and appearance easily. The main idea is then to *learn a good mapping network* (like current methods for 2D GAN inversion [44,49]) *that can encode the semantic mask into the somewhat smaller latent space domain for 3D controllable translation and manipulation*. As for the second issue, we intriguingly discover that a *region-aware learning strategy* is of vital importance. We therefore aim to tame an encoder that is sensitive to image patches, and adopt a region-based sampling pattern for the decoder. Furthermore, augmenting the input semantic masks with extracted contours and distance field representations [6] also considerably helps to highlight the intended semantic changes, making them more easily perceptible.

Following the above analysis, we build our Sem2NeRF framework upon the Swin Transformer encoder [33] and the pre-trained π -GAN decoder [3]. To kick off the single-view Semantic-to-NeRF translation task, we pinpoint two suitable yet challenging datasets, including CelebAMask-HQ [27] and CatMask, where the latter contains cat faces rendered using π -GAN and labelled with 6-class semantic masks using DatasetGAN [65]. We showcase the superiority of our model over several strong baselines by considering SofGAN [4], pix2pixHD [53] with GAN-inversion [25], and pSp [44]. Our contributions are three-fold:

- We introduce a novel and challenging task, Semantic-to-NeRF translation, which converts a single-view 2D semantic mask to a 3D scene modelled by neural radiance fields.
- With the insight of needing a region-aware learning strategy, we propose a novel framework, Sem2NeRF, which is capable of achieving 3D-consistent free viewpoint image generation, semantic editing and multi-model synthesis, by taking as input only one single-view semantic mask of a specific category, *e.g.*, human face, cat face.
- We validate our insight regarding our region-aware learning strategy and the efficacy of Sem2NeRF via extensive ablation studies, and demonstrate that Sem2NeRF outperforms strong baselines on two challenging datasets.

2 Related Work

NeRF and Generative NeRF. Starting as an approach focused on modelling a single static scene, NeRF [38] had seen rapid development in different aspects. Several approaches managed to reduce the training [50] and inference

time [31,35], while others improved visual quality [1]. Besides, it had also been extended in other ways, *e.g.*, dynamic scene [43], compositional scene [55], pose estimation [60], portrait generation [32], semantic segmentation [68].

Follow-up works that integrated NeRF with generative models were most relevant to ours. Schwarz *et al.* [46] proposed to learn a NeRF distribution by conditioning the input point positions with a sampled random vector. Niemeyer *et al.* [39] enabled multi-object generation by representing the whole scenes as a composition of different components. To improve the visual quality, π -GAN [3] adopted a SIREN-based [48] network structure with FiLM [42] conditioning. StyleNeRF [15] turned to embedding the volume rendering technique into StyleGAN [25]. More recently, VolumeGAN [59] relied on separately learning structure and texture features. MVCGAN [64] leveraged the underlying 3D geometry information. EG3D [2] proposed an efficient tri-plane hybrid 3D representation.

Our work belongs to the class of generative models, but unlike all existing methods that aimed to create a *random* scene, we aim to generate a *specific* scene that is conditioned by a given single-view semantic mask. Although there are concurrent works, *e.g.*, 3D-SGAN [61], FENeRF [51], exploring the similar condition settings, most of them purely focus on improving the quality of the generated images, while resort to existing GAN inversion [25] to do the mapping. In contrast, our work is more focused on improving the mapping from the mask to the NeRF-based scene.

Image-to-Image Translation is about converting an image from one source representation, *e.g.*, semantic masks, to another target representation, *e.g.*, photorealistic images. Since its introduction [20], progress has been made with regard to better image quality [53,7], multi-modal outputs [71,66,9], unsupervised learning [70,30], *etc.*. More recently, there is a new trend [44,47,58] of tackling this task by editing the latent space of a pre-trained generator, *e.g.*, StyleGAN.

In contrast to all mentioned work that aimed to map a semantic mask to an image, ours is focused on mapping to a 3D scene. We also notice that there are some recent approaches targeted at converting semantic masks to 3D scenes. Huang *et al.* [19] introduced rendering novel-view photorealistic images from a given semantic mask, by first applying semantic-to-image translation [40], then converting the single-view image to a 3D scene modelled by multiplane images (MPI) [69]. Hao *et al.* [16] proposed to learn a mapping from a semantically-labelled 3D block world to a NeRF-based 3D scene, using a scene-specific setting. Chen *et al.* [4] introduced a 3D-aware portrait generator by first mapping the given latent code to a semantic occupancy field (SOF) [8] for rendering novel view semantic masks, followed by applying image-to-image translation.

Unlike all mentioned attempts on learning semantic to 3D scene mappings, ours is the first to introduce the single-view semantic to NeRF translation task. Our work differs from theirs in: 1) We do not rely on any separate image-to-image translation stage, resulting in better multi-view consistency; 2) We do not require multi-view semantic masks for both training and testing phases, easing the data collection effort; 3) We pinpoint a solution for creating pseudo labels and demonstrate reasonable results beyond the human face domain.

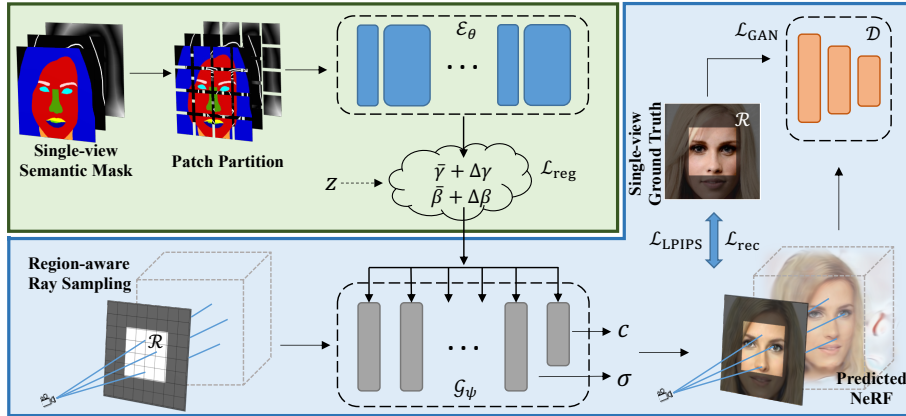


Fig. 2. Architecture of the Sem2NeRF framework. It aims to convert a single-view semantic mask to a 3D scene represented by NeRF. Specifically, a given semantic mask will be partitioned into patches, which will be further encoded by a patch-based encoder \mathcal{E}_θ into a latent style code (γ, β) of a pre-trained NeRF-based 3D generator \mathcal{G}_ψ . A region \mathcal{R} will be randomly sampled to enforce awareness of differences among regions. And an optional latent vector \mathbf{z} is included to enable multi-modal synthesis

3 Methodology

As shown in Fig. 1, our main goal is to train a Semantic-to-NeRF translation network $\Phi_{s \rightarrow \mathcal{V}}$, such that when presented with a single-view 2D semantic mask \mathbf{s} , it generates the corresponding NeRF representation \mathcal{V} , which can then be used to render realistic 3D-consistent images. This task is conceptually similar to the conventional semantic-to-image setting, except that here we opt to go beyond 2D image translation, and deal with the novel *controllable 3D translation*. More importantly, we can freely change the 3D content by *simply modifying the corresponding content in a single-view 2D semantic mask*.

In order to learn such a framework without enough supervision for arbitrary view appearances, we observed that 3D information can be learned from large image collections [22,3,2]. Therefore, our *key motivational insight* is this: instead of directly training $\Phi_{s \rightarrow \mathcal{V}}$ using *single-view* semantic-image pairs (\mathbf{s}, \mathbf{I}) (like current methods for 2D semantic-to-image translation [53,40]), we will train it as a two-stage pipeline shown in Fig. 2. Here, (A) we utilize a pre-trained 3D generator (lower portion \mathcal{G}_ψ) that learns 3D shape and appearance information from a large set of collected images; (B) we pose this challenging task as a *3D inversion* problem, where our main target is to design a front-end encoder (upper portion \mathcal{E}_θ) that maps the semantic mask into the generator latent space accurately.

The two training stages are executed independently and can be separately implemented with different frameworks. There are at least two unique benefits of breaking down the entire controllable 3D translation into two-stages: 1) The training does *not* require copious views of semantic-image pairs for each instance,

which are difficult to collect, or even impossible in some scenarios; 2) The compartmentalization of the 3D generator and the 2D encoder allows greater agility, where the 3D information can be previously learned on various tasks with a large collection of images and then be freely plugged into the 3D inversion pipeline.

3.1 3D Decoder with Region-Aware Ray Sampling

Preliminaries on NeRF. We first provide some preliminaries on NeRF before discussing how we exploit it for Sem2NeRF. NeRF [38] is one kind of implicit functions that represents a continuous 3D scene, which has achieved great successes in modeling 3D shape and appearance. A NeRF is a neural network that maps a 3D location $\mathbf{x} \in \mathbb{R}^3$ and a viewing direction $\mathbf{d} \in \mathbb{S}^2$ to a spatially varying volume density σ and a view-dependent emitted color $\mathbf{c} = (r, g, b)$. NeRFs trained on natural images are able to continuously render realistic images at arbitrary views. In particular, it requires to use the volume rendering [28], which computes the following integral to obtain the color of a pixel:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds\right), \quad (1)$$

where $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is the ray casting from the virtual camera located at \mathbf{o} , bounded by near t_n and far t_f , and $T(t)$ represents the accumulated transmittance of the ray traveling from t_n to t . The integral $C(\mathbf{r})$ is further implemented with a hierarchical volume sampling strategy [28,38], resulting in the optimization of a ‘‘coarse’’ network followed by a ‘‘fine’’ network.

NeRF-based Generator. Our work is mainly based on a representative NeRF-based generator, π -GAN [3], which learns 3D representation using only 2D supervision. Inspired by StyleGAN2 [25], the architecture of π -GAN is mainly composed of two parts, a mapping network $\mathcal{F} : \mathcal{Z} \rightarrow \mathcal{W}$ that maps a latent vector z in the input latent space \mathcal{Z} to an intermediate latent vector $w \rightarrow \mathcal{W}$, and a SIREN-based [48] synthesis network that maps w to the NeRF representation \mathcal{V} that supports rendering 3D-consistent images from arbitrary camera poses.

Our Sem2NeRF framework can use various NeRF-based generators. Here, we choose π -GAN as the main decoder in our architecture for two main reasons. Firstly, among all *published* works related to NeRF-based generators, π -GAN achieves state-of-the-art performance in terms of rendered image quality and their underlying 3D consistency. Secondly and more importantly, similar to StyleGAN, the FiLM [42] conditioning used by π -GAN enables layer-wise control over the decoder and the mapping network decouples some high-level attributes, making it easier to perform *3D inversion* on top of NeRF, *i.e.*, searching for the desired latent code w that best reconstructs an ideal target. The similar observation has been previously explored in the latest 2D GAN inversion [10,25,44].

Region-Aware Ray Sampling. While π -GAN already provides high-quality view-consistent rendered images, our main goal is to accurately restore the NeRF

from a single-view semantic mask, and even freely edit the 3D content via such a map. To achieve this, *the network should be sensitive to local small modifications*. However, this is not supported in the original π -GAN, which is trained on each entire image with a global perception. It stores scene-specific information in a latent code, which is shared across all points that are bounded by the rendering volume. As a result, a small change in an original latent code will easily cause a global modification in generation. This may not impact pure 3D generation, for which only the quality of global shape and appearance is paramount, but it has a large negative effect on recreating a 3D representation that accurately matches the corresponding semantic mask.

To mitigate this issue, we adopt a region-based ray sampling pattern [46,32] in the π -GAN decoder, that *attempts to encourage latent codes to represent local regions at different scales and locations*. Suppose the rendered image \mathbf{I} with a target size $h \times w$, a local region \mathcal{R} used for training is randomly sampled as

$$\mathcal{R}(\alpha, (\Delta h, \Delta w)) = \{(\alpha h + \Delta h, \alpha w + \Delta w)\}, \quad (2)$$

where $(\alpha h + \Delta h, \alpha w + \Delta w)$ denotes the sampling coordinates of rays, with $\alpha \in (0, 1]$ being the scaling factor and $(\Delta h \in [0, (1 - \alpha)h], \Delta w \in [0, (1 - \alpha)w])$ being the translation factor. To obtain such training pairs between the NeRF rendered output and the local ground truth, we sample the original whole image using the same region coordinates \mathcal{R} with bilinear interpolation. This strategy leads to large improvements on conditional generation as shown in the experiments.

3.2 3D Inversion Using Region-Aware 2D Encoder

3D Inversion. To inversely map a semantic mask \mathbf{s} into the \mathcal{W} latent space of the 3D generator \mathcal{G}_ψ by an encoder \mathcal{E}_θ , with respective parameters ψ and θ , we train \mathcal{E}_θ to minimize the reconstruction error between ground truth image \mathbf{I} and output $\hat{\mathbf{I}}$. Specifically, Semantic-to-NeRF translation represents the mapping

$$\Phi_{\mathbf{s} \rightarrow \mathcal{V}}(\mathbf{x}, \mathbf{d}, \mathbf{z}; \mathbf{s}) = \mathcal{G}_\psi(\mathbf{x}, \mathbf{d}, \mathbf{z}; \mathcal{E}_\theta(\mathbf{s})) = \mathcal{V}(\sigma, \mathbf{c}) \quad (3)$$

where \mathbf{x}, \mathbf{d} denotes point position and ray direction, while the derived density σ and color \mathbf{c} can be used to calculate the corresponding pixel value via volume rendering as in Eq. (1). For *controllable* 3D generation, \mathbf{s} is the input single-view semantic mask, embedded into \mathcal{W} space to control the generated 3D content, while we also enable multi-modal synthesis by adding another latent vector \mathbf{z} to model the generated appearance. Note that \mathbf{s} only comes in a single view, which is not necessary the same as the output viewing direction. In short, *we use only single-view semantic-image pairs (\mathbf{s}, \mathbf{I}) for the Sem2NeRF training*, as the 3D view-consistent information has been captured by the *fixed* pre-trained 3D generator \mathcal{G}_ψ . Hence, we focus only on training the encoder network \mathcal{E}_θ to learn the posterior distribution $q(w|\mathbf{s})$ for 3D inversion.

Region-Aware 2D Encoder. A simple way for 3D inversion is to directly apply an existing 2D GAN inversion framework. However, this straightforward idea

does *not* work well as we originally discovered when using the state-of-the-art pSp encoder [44] in our setting, especially for small but perceptually important regions, such as eyes. Our conjecture is that the conventional CNN-based architecture integrates the neighboring information via overaggressive filtering, resulting in heavy loss of small details [62].

To mitigate this issue, we also deploy a region-aware learning strategy in the 2D encoder, which is inspired by the latest patch-based methods [12,67] that capture information in every patch with equal possibility. In other words, when we directly extract features from local patches, it will be *more sensitive to the semantic variation within each patch*, which can ameliorate the problem of imbalanced semantic distribution within an image. In particular, we adopt the Swin Transformer [33] as the encoder architecture. To embed the semantic mask \mathbf{s} into the \mathcal{W} latent space of the pre-trained 3D generator, we replace the final classification output size with the size of the latent vectors \mathbf{w} . Besides, to further stabilize the inversion training, we take inspiration from the truncation trick [24,44] and set the learned latent codes for the pre-trained decoder as

$$\gamma = \bar{\gamma} + \Delta\gamma, \beta = \bar{\beta} + \Delta\beta, \quad (4)$$

where γ and β represent the embedded vectors for the \mathcal{W} latent space, *i.e.*, frequency and phase shift of π -GAN, respectively; $\Delta\gamma$ and $\Delta\beta$ are the outputs of the proposed encoder \mathcal{E}_θ , while $\bar{\gamma}$ and $\bar{\beta}$ are the average latent codes extracted by the pre-trained π -GAN original mapping network $\mathcal{F} : \mathcal{Z} \rightarrow \mathcal{W}$.

Additional Inputs for the 2D Encoder. As mentioned, a semantic mask contains sparse information, where the changing of small regions may be imperceptible to the network, making the semantic-based controllable 3D editing very challenging. Considering that editing a semantic mask only effectively alters the boundaries between different semantic labels, we conjecture that explicitly augmenting the semantic input with *boundary information* will be useful for semantic editing. Therefore, we concatenate the semantic mask input with contours and distance field representations [6] for the region-aware encoder. These additional inputs further improve the semantic editing performance considerably as shown in the experiments. Note that contours and distance field representations are both directly calculated from the semantic masks (refer to Section A.1 for more details), which do *not involve any extra labels*.

3.3 Training Loss Functions

During the training phase, we use the single-view semantic mask \mathbf{s} , the corresponding viewing direction d_s , and the paired ground truth RGB image \mathbf{I} . Similar to Semantic-to-Image translation, we start by applying a pixel-level reconstruction loss,

$$\mathcal{L}_{\text{rec}}(\mathbf{I}, \mathbf{s}, d_s) = \|\mathbf{I} - \mathcal{G}_\psi(\mathcal{E}_\theta(\mathbf{s}), d_s)\|_2, \quad (5)$$

where $\mathcal{E}_\theta(\mathbf{s})$ denotes the latent codes mapped from \mathbf{s} via the region-aware encoder $\mathcal{E}_\theta(\cdot)$, while $\mathcal{G}_\psi(\mathcal{E}_\theta(\mathbf{s}), d_s)$ represents the generated image rendered from direction

d_s via the decoder $\mathcal{G}_\psi(\cdot)$. Unless otherwise specified, the aforementioned region-aware sampling strategy is applied to \mathcal{G}_ψ and \mathbf{I} before calculating any losses.

To further enforce the feature-level similarity between the generated image and the ground truth, the LPIPS loss [63] is leveraged,

$$\mathcal{L}_{\text{LPIPS}}(\mathbf{I}, \mathbf{s}, d_s) = \|\mathcal{F}(\mathbf{I}) - \mathcal{F}(\mathcal{G}_\psi(\mathcal{E}_\theta(\mathbf{s}), d_s))\|_2, \quad (6)$$

where $\mathcal{F}(\cdot)$ refers to the pre-trained feature extraction network.

Inspired by the truncation trick [24,44], we further encourage the decoder latent codes γ, β to be close to the average codes $\bar{\gamma}, \bar{\beta}$, which is achieved by regularizing the encoder with

$$\mathcal{L}_{\text{reg}}(\mathbf{s}) = \|\mathcal{E}_\theta(\mathbf{s})\|_2. \quad (7)$$

To improve image quality, especially for novel views, we further apply a non-saturating GAN loss with R1 regularization [37],

$$\mathcal{L}_{\text{GAN}}(\mathbf{I}, \mathbf{s}, d) = f(\mathcal{D}(\mathcal{G}_\psi(\mathcal{E}_\theta(\mathbf{I}), d))) + f(-\mathcal{D}(\mathbf{I})) + \lambda_{\text{R1}} |\nabla \mathcal{D}(\mathbf{I})|^2, \quad (8)$$

where $f(u) = -\log(1 + \exp(-u))$.

Here $\mathcal{D}(\cdot)$ is a patch discriminator [20], aligned with our region-aware learning strategy for the decoder, and λ_{R1} is a hyperparameter that is set to 10. Note that here the viewing direction d is not required to be the same as the input semantic viewing direction d_s , and we randomly sample this viewing direction from a known distribution, *i.e.* Gaussian, following the settings of π -GAN [3].

Finally, the overall training objective for our framework is a weighted combination of the above loss functions as

$$\mathcal{L}_{\text{Sem2NeRF}} = \lambda_{\text{rec}} \mathcal{L}_{\text{rec}} + \lambda_{\text{LPIPS}} \mathcal{L}_{\text{LPIPS}} + \lambda_{\text{reg}} \mathcal{L}_{\text{reg}} + \lambda_{\text{GAN}} \mathcal{L}_{\text{GAN}}. \quad (9)$$

3.4 Model Inference

For inference, our model takes as input a 2D single-view semantic mask, while d_s is optional, required only when rendering an image with the same viewing direction as the semantic mask. Different from the training phase, during inference the rays are cast to cover the whole image plane, rather than a local region.

Multi-View Generation. Since the employed decoder is a NeRF-based generator, Sem2NeRF inherently supports novel view generation. Specifically, given a semantic mask \mathbf{s} , it will first be mapped as an embedded vector in the \mathcal{W} latent space that controls the “content” of the NeRF-based generator, whereupon a novel view image can then be generated by volume rendering the NeRF from an arbitrary viewing direction.

Multi-Modal Synthesis. Similar to the diversified mapping in semantic-to-image [40], ideally a single semantic mask should be translated into multiple NeRFs consistent to it. Our Sem2NeRF framework inherently supports multi-modal synthesis in inference due to the usage of FiLM [42] conditioning on π -GAN, without requiring any special customization in training. In practice, we additionally pass a random-sampled vector to the pre-trained π -GAN noise mapping module to obtain corresponding latent style codes \mathbf{z} . Style mixing [44,24] is then performed between \mathbf{z} and $\mathcal{E}_\theta(\mathbf{s})$ to yield multi-modal outcomes.

4 Experiments

4.1 Settings

Datasets. To achieve Semantic-to-NeRF translation, we assume the training data to have single-view registered semantic masks and images, with the corresponding viewing directions. Two datasets were used for evaluation in our experiments. **CelebAMask-HQ** [27] contains images from CelebA-HQ [34,23], manually-labelled 19-class semantic masks, and head poses. We merged the left-right labels of symmetric parts, *i.e.*, eyes, eyebrows and ears, into one label per part. The dataset was randomly partitioned into training set with 28,000 samples and test set with 2,000 samples. **CatMask** is built using π -GAN and Dataset-GAN [65] to further demonstrate the potential of Semantic-to-NeRF task and Sem2NeRF. Technical details are elaborated in Section A.2.

Baselines. We identified the following three methods as baselines for comparison in our introduced Semantic-to-NeRF task. **SofGAN** [4] is an image translation approach. For a given single-view mask, we first apply inversion via iterative optimizations to find the corresponding latent vector for the preceding SOF [8] network, which can generate novel view semantic masks for further image-to-image mapping. Note that SofGAN requires training data to have high-quality multi-view semantic masks, which is not available nor needed in our task. **pix2pixHD** [53] is an image translation approach. We adopt it with general GAN-inversion techniques [3,25]. For a given mask, it is first mapped to a photo-realistic image via pix2pixHD, which will then be mapped to the corresponding latent code in π -GAN via GAN-inversion. With the recovered codes, multi-view images can be directly obtained using π -GAN. **pSp** [44] is an image translation approach that is designed for encoding into StyleGAN2 [25]. We adapted it by using its ResNet [17]-based pSp encoder to replace the π -GAN mapping network, and we further trained the network with objective functions used by pSp.

Evaluation Metrics. We show qualitative results by rendering images with different viewing directions and FOV (Field of View). We also report Frechet Inception Distance (FID) [18] and Inception Score (IS) [45] using Inception-v3 [52] over the test sets. Average running time and model sizes are also compared.

Implementation Details. Swin-T is used in all experiments with input resolution 224×224 . For the decoder, the size of local region \mathcal{R} is set to 128×128 . The step size of each ray is set to 28. Other miscellaneous settings of the pre-trained decoder, *e.g.*, ray depth ranges, are kept unchanged. Hyper-parameters in Eq. (9) are set as $\lambda_{\text{rec}}=1$, $\lambda_{\text{LPIPS}}=0.8$, $\lambda_{\text{reg}}=0.005$, $\lambda_{\text{GAN}}=0.08$. The implementation is done in PyTorch [41]. More details are provided in Section A.3.

4.2 Results

Comparisons on CelebAMask-HQ. As shown in Fig. 3, compared to all other baseline models, **Sem2NeRF (1st&5th columns)** achieved the best performance on both mapping accuracy and multi-view consistency. **pSp (2nd&6th columns)** generated images with lower quality compared to ours, especially for

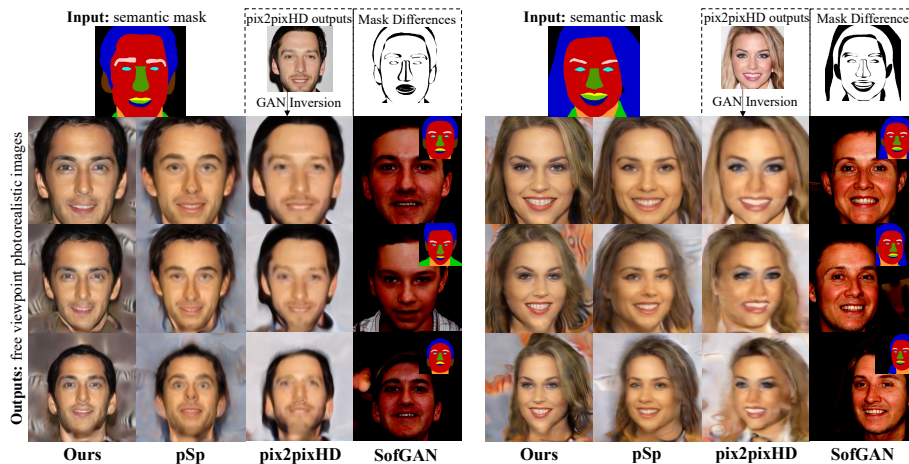


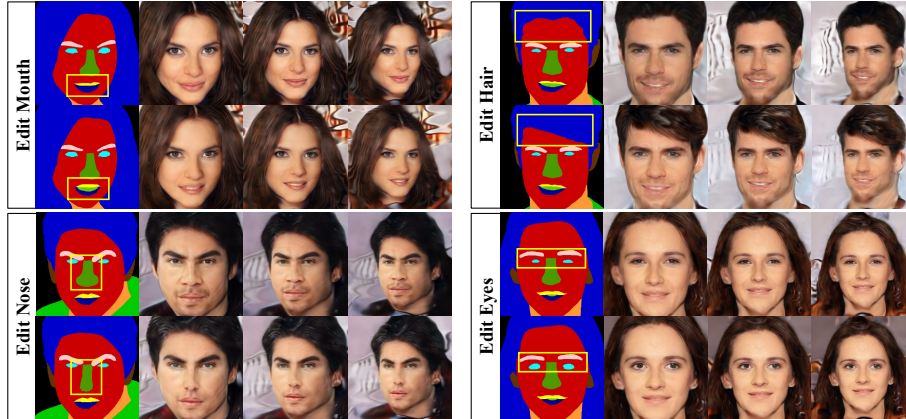
Fig. 3. Comparisons on CelebAMask. Images at each column are generated by the corresponding models mentioned at the bottom. Only SofGAN requires generation of multi-view semantic masks, shown at the top right corners of related images

novel views, mainly because our model is designed with a region-aware learning strategy and a GAN loss for random-posed images during training. The CNN-based encoder also failed to capture fine-grained details, *e.g.*, eyebrow shapes for the left face. Our method and the inversion-based pix2pixHD were better in matching semantics compared to pSp. **pix2pixHD (3rd&7th columns)** can map semantic masks to high quality images in the same viewpoint (top row), but does not generate novel views well. Basic GAN-inversion is not an efficient or easy way to find the desired latent codes, since the current 3D generative models are still quite immature. Even though images with the same viewing direction as the masks are reasonable, those novel view outputs contain artifacts. **SofGAN (4th&8th columns)** generated each single-view image with good quality; however, its results do not match with the given mask and lacked 3D consistency. The reason is that it is hard to map the given semantic mask to the desired latent codes of its semantic generator (SOF Net), whose sampling space is relatively small due to the lack of training data (only 122 subjects). The recovered mask did not match well with the given mask (top row). Besides, although the semantic masks show good multi-view consistency (top right corner of each image), conducting semantic-to-image mapping separately for each viewpoint does not guarantee that the consistency will be retained, since a semantic mask hardly contains any texture information and is geometrically ambiguous.

Quantitative results are give in Table 1. It can be seen that our Sem2NeRF method achieves the best performance, significantly outperforming the two base-lines in both FID and IS scores. Note that we did not quantitatively compare single view image quality with SofGAN, considering that SofGAN for Semantic-to-NeRF is limited by its mask inversion quality and multi-view consistency,

Table 1. Quantitative comparisons on CelebAMask @ 128×128

	FID ↓	IS ↑	Runtime(s) ↓	# Params(M) ↓
pix2pixHD [53] (with inversion)	67.32	1.72	161.59±0.859	~184.24
pSp [44]	55.56	1.74	0.25±0.004	~138.27
Sem2NeRF(Ours)	41.52	2.03	0.18±0.003	~32.01

**Fig. 4.** Editing 3D scenes by changing single-view semantic masks. Three viewpoints are shown for better comparison in each group

both of which cannot be measured by FID or IS scores. We also notice that scores of all models are lower than expected. The main reason is that π -GAN is initially trained on CelebA, but due to the requirement of semantic masks, our task conducted experiments using CelebA-HQ. The domain gap between CelebA and CelebA-HQ reduced the FID scores dramatically. Besides, our model also sees advantages in terms of running time and model size.

Mask Editing. As depicted in Fig. 4, our framework supports editing of 3D scenes by simply changing the given semantic mask, and is applicable to both labels associated with large regions, *e.g.*, hair, as well as small regions, *e.g.*, eyes, nose, mouth. This is not trivial since the semantic mask is not directly leveraged to control the 3D scene at the pixel level (if even possible), but is instead encoded into a sparse latent code, which may fail to preserve fine-grain editing. We address this challenge via the region-aware learning strategy.

Multi-Modal Synthesis. Sem2NeRF supports multi-modal synthesis by simply changing the last few layers of the style codes. As shown in Fig. 5, we randomly sampled two style codes, and applied linear blending to continuously change the general styles of the 3D scenes generated by the given masks.

Ablation Studies. To further evaluate the efficacy of Sem2NeRF, we designed four ablation models, including 1) without region-aware encoder $\Phi_{\text{wo_RE}}$, where the Swin-T encoder is replaced by the pSp encoder; 2) without region-aware



Fig. 5. Multi-modal synthesis. Styles are linearly blended from left to right. Three viewpoints are provided from top to bottom

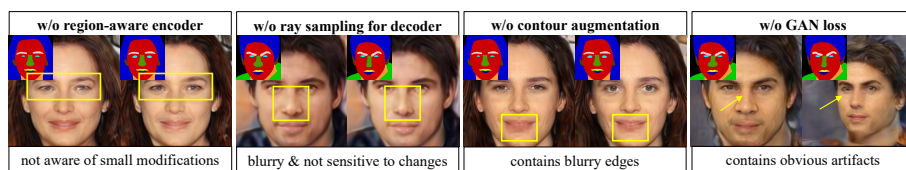


Fig. 6. Results of ablation studies. Each group (two views) is generated by a model without the component mentioned at the top. Main issues are described at the bottom

decoder $\Phi_{\text{wo_RD}}$, where the region-aware ray sampling strategy is discarded; 3) without input augmentation $\Phi_{\text{wo_IA}}$, where contours and distance field representations are removed from the input; and 4) without random-pose GAN loss $\Phi_{\text{wo_GAN}}$, where both Eq. (8) and the discriminator are removed.

As shown in Fig. 6, compared to the full model (Fig. 4), $\Phi_{\text{wo_RE}}$ (1st group) is not sensitive to changes in small regions, *i.e.*, eyes, mainly because the CNN-based encoder tends to ignore small changes. $\Phi_{\text{wo_RD}}$ (2nd group) shows similar pattern (nose region) as the latent codes are not trained to be region-aware. It also has lower image quality, because the region-aware strategy enables denser sampling. $\Phi_{\text{wo_IA}}$ (3rd group) achieves comparable performance but with blurry edges for some regions, *e.g.*, mouth. This is because both contour and distance field representation help highlight the boundary information. Finally, images obtained by $\Phi_{\text{wo_GAN}}$ (4th group) have more artifacts in both views, demonstrating that GAN loss is important for improving the image quality of different poses.

Experiments beyond Human Faces. The introduced task can easily go beyond the human face domain by leveraging state-of-the-art weakly supervised semantic segmentation model to create pseudo labels. In this work, we present a Cat face example. Experimental results are shown in Fig. 7. Even when training with noisy pseudo labels, Sem2NeRF is robust enough to generate plausible results. For a given cat semantic mask, our model can map it to a 3D scene and render cat faces from arbitrary viewpoints, including different viewing directions



Fig. 7. Results on CatMask. Left part compares results of changing eyes shape. Right part showcases results of style linear blending (in zigzag order)

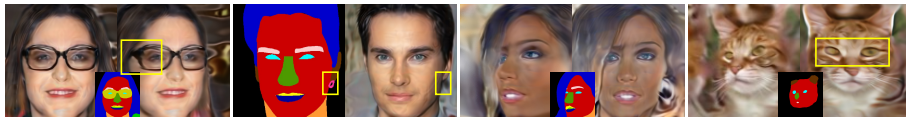


Fig. 8. Challenging cases of Sem2NeRF on Semantic-to-NeRF translation

(left part), and different FOV (right part). It also allows changing the 3D scenes by editing the single-view semantic masks, *e.g.*, changing the eye shape (left two rows). Multi-modal synthesis is also supported (right part in zigzag order).

Challenging Cases. Although Sem2NeRF addresses the Semantic-to-NeRF task in most cases, its advantages rely on an assumption, namely the generative capability of the pre-trained decoder. We show some challenging cases in Fig. 8. Accessories may have the wrong geometric shape (glasses in 1st case), or fail to render (earring in 2nd case), while masks with extreme poses might be converted to 3D scenes with abnormal texture or distorted contents (last two cases).

5 Conclusions

We have presented an initial step of extending the 2D image-to-image task to the 3D space, and introduced a new task called Semantic-to-NeRF translation. It aims to reconstruct a NeRF-based 3D scene, by taking as input only one single-view semantic mask. We further proposed Sem2NeRF model, which addresses the task via encoding the semantic mask into the latent space of a pre-trained 3D generative model. More importantly, we intriguingly found the importance of regional awareness for this new task, and tamed Sem2NeRF with a region-aware learning strategy. We demonstrated the capability of Sem2NeRF regarding free viewpoint generation, mask editing and multi-modal synthesis on two benchmark datasets, and showcased the superiority of our framework compared to three strong baselines. Future work will include adding more scenarios to the new task, and supporting changing styles for specific regions.

Acknowledgements This research is partially supported by CRC Building 4.0 Project #44.

References

1. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In: *Int. Conf. Comput. Vis.* pp. 5855–5864 (2021)
2. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., De Mello, S., Gallo, O., Guibas, L., Tremblay, J., Khamis, S., et al.: Efficient geometry-aware 3d generative adversarial networks. *IEEE Conf. Comput. Vis. Pattern Recog.* (2022)
3. Chan, E.R., Monteiro, M., Kellnhofer, P., Wu, J., Wetzstein, G.: pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 5799–5809 (2021)
4. Chen, A., Liu, R., Xie, L., Chen, Z., Su, H., Yu, J.: Sofgan: A portrait image generator with dynamic styling. *ACM Trans. Graph.* **41**(1), 1–26 (2022)
5. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Trans. Pattern Anal. Mach. Intell.* **40**(4), 834–848 (2017)
6. Chen, W., Hays, J.: Sketchygan: Towards diverse and realistic sketch to image synthesis. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 9416–9425 (2018)
7. Chen, Y., Huang, J., Wang, J., Xie, X.: Edge prior augmented networks for motion deblurring on naturally blurry images. *arXiv preprint arXiv:2109.08915* (2021)
8. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 5939–5948 (2019)
9. Choi, Y., Uh, Y., Yoo, J., Ha, J.W.: Stargan v2: Diverse image synthesis for multiple domains. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 8188–8197 (2020)
10. Collins, E., Bala, R., Price, B., Susstrunk, S.: Editing in style: Uncovering the local semantics of gans. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 5771–5780 (2020)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 248–255. *Ieee* (2009)
12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *Int. Conf. Learn. Represent.* (2021)
13. Goel, S., Kanazawa, A., Malik, J.: Shape and viewpoint without keypoints. In: *Eur. Conf. Comput. Vis.* pp. 88–104. Springer (2020)
14. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. *Adv. Neural Inform. Process. Syst.* **27** (2014)
15. Gu, J., Liu, L., Wang, P., Theobalt, C.: Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *Int. Conf. Learn. Represent.* (2022)
16. Hao, Z., Mallya, A., Belongie, S., Liu, M.Y.: Gancraft: Unsupervised 3d neural rendering of minecraft worlds. In: *Int. Conf. Comput. Vis.* pp. 14072–14082 (2021)
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 770–778 (2016)
18. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Adv. Neural Inform. Process. Syst.* **30** (2017)
19. Huang, H.P., Tseng, H.Y., Lee, H.Y., Huang, J.B.: Semantic view synthesis. In: *Eur. Conf. Comput. Vis.* pp. 592–608. Springer (2020)

20. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 1125–1134 (2017)
21. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: *Eur. Conf. Comput. Vis.* pp. 694–711. Springer (2016)
22. Kanazawa, A., Tulsiani, S., Efros, A.A., Malik, J.: Learning category-specific mesh reconstruction from image collections. In: *Eur. Conf. Comput. Vis.* pp. 371–386 (2018)
23. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. *Int. Conf. Learn. Represent.* (2018)
24. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 4401–4410 (2019)
25. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 8110–8119 (2020)
26. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *Int. Conf. Learn. Represent.* (2014)
27. Lee, C.H., Liu, Z., Wu, L., Luo, P.: Maskgan: Towards diverse and interactive facial image manipulation. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 5549–5558 (2020)
28. Levoy, M.: Efficient ray tracing of volume data. *ACM Trans. Graph.* **9**(3), 245–261 (1990)
29. Ling, H., Kreis, K., Li, D., Kim, S.W., Torralba, A., Fidler, S.: Editgan: High-precision semantic image editing. In: *Adv. Neural Inform. Process. Syst.* (2021)
30. Lira, W., Merz, J., Ritchie, D., Cohen-Or, D., Zhang, H.: Ganhopper: Multi-hop gan for unsupervised image-to-image translation. In: *Eur. Conf. Comput. Vis.* pp. 363–379. Springer (2020)
31. Liu, L., Gu, J., Zaw Lin, K., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. *Adv. Neural Inform. Process. Syst.* **33**, 15651–15663 (2020)
32. Liu, X., Xu, Y., Wu, Q., Zhou, H., Wu, W., Zhou, B.: Semantic-aware implicit neural audio-driven video portrait generation. *Eur. Conf. Comput. Vis.* (2022)
33. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: *Int. Conf. Comput. Vis.* pp. 10012–10022 (2021)
34. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: *Int. Conf. Comput. Vis.* pp. 3730–3738 (2015)
35. Lombardi, S., Simon, T., Schwartz, G., Zollhoefer, M., Sheikh, Y., Saragih, J.: Mixture of volumetric primitives for efficient neural rendering. *ACM Trans. Graph.* **40**(4), 1–13 (2021)
36. Luo, S., Hu, W.: Diffusion probabilistic models for 3d point cloud generation. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 2837–2845 (2021)
37. Mescheder, L., Geiger, A., Nowozin, S.: Which training methods for gans do actually converge? In: *Int. Conf. Mach. Learn.* pp. 3481–3490. PMLR (2018)
38. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: *Eur. Conf. Comput. Vis.* pp. 405–421. Springer (2020)
39. Niemeyer, M., Geiger, A.: Giraffe: Representing scenes as compositional generative neural feature fields. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 11453–11464 (2021)

40. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 2337–2346 (2019)
41. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inform. Process. Syst.* **32** (2019)
42. Perez, E., Strub, F., De Vries, H., Dumoulin, V., Courville, A.: Film: Visual reasoning with a general conditioning layer. In: *AAAI*. vol. 32 (2018)
43. Pumarola, A., Corona, E., Pons-Moll, G., Moreno-Noguer, F.: D-nerf: Neural radiance fields for dynamic scenes. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 10318–10327 (2021)
44. Richardson, E., Alaluf, Y., Patashnik, O., Nitzan, Y., Azar, Y., Shapiro, S., Cohen-Or, D.: Encoding in style: a stylegan encoder for image-to-image translation. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 2287–2296 (2021)
45. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. *Adv. Neural Inform. Process. Syst.* **29** (2016)
46. Schwarzer, K., Liao, Y., Niemeyer, M., Geiger, A.: Graf: Generative radiance fields for 3d-aware image synthesis. *Adv. Neural Inform. Process. Syst.* **33**, 20154–20166 (2020)
47. Shi, Y., Yang, X., Wan, Y., Shen, X.: Semanticstylegan: Learning compositional generative priors for controllable image synthesis and editing. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 11254–11264 (2022)
48. Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.: Implicit neural representations with periodic activation functions. *Adv. Neural Inform. Process. Syst.* **33**, 7462–7473 (2020)
49. Song, G., Luo, L., Liu, J., Ma, W.C., Lai, C., Zheng, C., Cham, T.J.: Agilegan: stylizing portraits by inversion-consistent transfer learning. *ACM Trans. Graph.* **40**(4), 1–13 (2021)
50. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. *IEEE Conf. Comput. Vis. Pattern Recog.* (2022)
51. Sun, J., Wang, X., Zhang, Y., Li, X., Zhang, Q., Liu, Y., Wang, J.: Fenerf: Face editing in neural radiance fields. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 7672–7682 (2022)
52. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 2818–2826 (2016)
53. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 8798–8807 (2018)
54. Wu, J., Zhang, C., Xue, T., Freeman, B., Tenenbaum, J.: Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Adv. Neural Inform. Process. Syst.* **29** (2016)
55. Wu, Q., Liu, X., Chen, Y., Li, K., Zheng, C., Cai, J., Zheng, J.: Object-compositional neural implicit surfaces. *Eur. Conf. Comput. Vis.* (2022)
56. Wu, Z., Lischinski, D., Shechtman, E.: Stylespace analysis: Disentangled controls for stylegan image generation. In: *IEEE Conf. Comput. Vis. Pattern Recog.* pp. 12863–12872 (2021)
57. Wu, Z., Nitzan, Y., Shechtman, E., Lischinski, D.: Stylealign: Analysis and applications of aligned stylegan models. *Int. Conf. Learn. Represent.* (2022)

58. Xu, Y., Yin, Y., Jiang, L., Wu, Q., Zheng, C., Loy, C.C., Dai, B., Wu, W.: Transeditor: Transformer-based dual-space gan for highly controllable facial editing. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 7683–7692 (2022)
59. Xu, Y., Peng, S., Yang, C., Shen, Y., Zhou, B.: 3d-aware image synthesis via learning structural and textural representations. IEEE Conf. Comput. Vis. Pattern Recog. (2022)
60. Yen-Chen, L., Florence, P., Barron, J.T., Rodriguez, A., Isola, P., Lin, T.Y.: inerf: Inverting neural radiance fields for pose estimation. In: IEEE Int. Conf. Intell. Robots Syst. pp. 1323–1330. IEEE (2021)
61. Zhang, J., Sangineto, E., Tang, H., Siarohin, A., Zhong, Z., Sebe, N., Wang, W.: 3d-aware semantic-guided generative model for human synthesis. Eur. Conf. Comput. Vis. (2022)
62. Zhang, R.: Making convolutional networks shift-invariant again. In: Int. Conf. Mach. Learn. pp. 7324–7334. PMLR (2019)
63. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 586–595 (2018)
64. Zhang, X., Zheng, Z., Gao, D., Zhang, B., Pan, P., Yang, Y.: Multi-view consistent generative adversarial networks for 3d-aware image synthesis. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 18450–18459 (2022)
65. Zhang, Y., Ling, H., Gao, J., Yin, K., Lafleche, J.F., Barriuso, A., Torralba, A., Fidler, S.: Datasetgan: Efficient labeled data factory with minimal human effort. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 10145–10155 (2021)
66. Zheng, C., Cham, T.J., Cai, J.: Pluralistic image completion. In: IEEE Conf. Comput. Vis. Pattern Recog. pp. 1438–1447 (2019)
67. Zheng, C., Cham, T.J., Cai, J.: Tfill: Image completion via a transformer-based architecture. IEEE Conf. Comput. Vis. Pattern Recog. (2022)
68. Zhi, S., Laidlow, T., Leutenegger, S., Davison, A.J.: In-place scene labelling and understanding with implicit scene representation. In: Int. Conf. Comput. Vis. pp. 15838–15847 (2021)
69. Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: Learning view synthesis using multiplane images. ACM Trans. Graph. (2018)
70. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Int. Conf. Comput. Vis. pp. 2223–2232 (2017)
71. Zhu, J.Y., Zhang, R., Pathak, D., Darrell, T., Efros, A.A., Wang, O., Shechtman, E.: Toward multimodal image-to-image translation. Adv. Neural Inform. Process. Syst. **30** (2017)

A Additional Technical Details

A.1 Builders for Additional Inputs

As mentioned in Section 3.2, additional inputs, *i.e.*, contour and distance field representation, are provided for the encoder of Sem2NeRF to further highlight the boundary information. In this subsection, we will detail how to build these data using the python package “cv2”³. Note that both of them are directly calculated from the semantic mask, without involving any extra labels. And each builder function can be done in 1 ~ 2 milliseconds.

Contour is generally represented as a curve, joining all the continuous points that share the same intensity. It is widely used as a tool to help shape analysis, object detection, *etc.*. In our implementation, we build the contour from the given one-hot encoded semantic mask. Main python codes are given as below. An output example is shown in Figure 9 (middle).

```

1 # ----- CONTOUR BUILDER -----
2 def binary_masks_to_contour(binary_masks):
3     ''' INPUT: semantic mask in one-hot encoding form
4         OUTPUT: a contour map of the given semantic mask '''
5
6     # initialize a black canvas
7     mask = numpy.zeros((512, 512, 3), dtype=numpy.uint8)
8     # find contours for each label
9     for binary_mask in binary_masks:
10        cnts = cv2.findContours(binary_mask,
11                               cv2.RETR_EXTERNAL,
12                               cv2.CHAIN_APPROX_SIMPLE)
13        cnts = cnts[0] if len(cnts) == 2 else cnts[1]
14        for c in cnts:
15            # draw contour with white color on the canvas
16            cv2.drawContours(mask, [c], -1, (255, 255, 255),
17                             thickness=3)
18        contour = cv2.cvtColor(mask, cv2.COLOR_BGR2GRAY)
19    return contour
20 # ----- END -----

```

Distance field representation is a dense representation extracted from binary image via distance transformation. In the distance field, the grey intensity of each pixel indicates its distance to the nearest boundary. An unsigned Euclidean distance field representation is adopted in our experiments. Main python codes are given as below. An output example is shown in Figure 9 (right).

```

1 # ----- DISTANCE FIELD REPRESENTATION BUILDER -----
2 def contour_to_dist_filed(contour):
3     ''' INPUT: contour, contour of the semantic mask
4         OUTPUT: dist_field, distance filed representation '''
5

```

³ <https://pypi.org/project/opencv-python/>

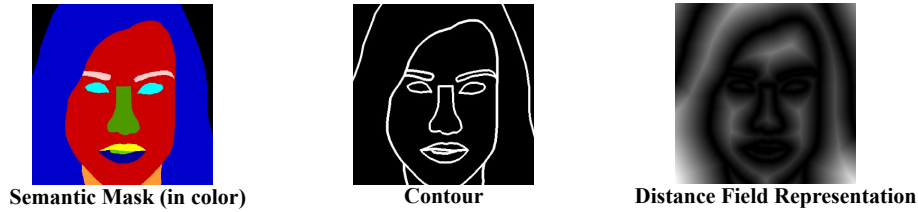


Fig. 9. An example of the encoder input. Semantic mask is shown in color for better visualization, while the network takes one-hot encoded mask as input

```

6   # invert background and foreground of contour to match
   the setting
7   invert_contour = cv2.bitwise_not(contour)
8   dist_field = cv2.distanceTransform(invert_contour,
9                                     cv2.DIST_L2, 3)
10  # normalize the distance field to [0, 1]
11  cv2.normalize(dist_field, dist_field, 0, 1.0,
12               cv2.NORM_MINMAX)
13  dist_field = dist_field * 255.
14  return dist_field
15 # ----- END -----

```

A.2 CatMask Dataset Rendering

To better demonstrate the introduced Semantic-to-NeRF translation task and evaluate the proposed Sem2NeRF framework, we develop a general solution to create datasets with pseudo labels using minimal human effort. In this work, we present an example by rendering a cat faces dataset, termed CatMask, which contains single-view cat faces generated by the pre-trained π -GAN model, pseudo ground-truth viewing directions, and 6-classes semantic masks labelled by DatasetGAN [65]. Similar to CelebAMask-HQ, CatMask only varies on the yaw and pitch axis. And it contains 28,000 training images and 2,000 testing images.

Cat faces from π -GAN. As mentioned in Section 4.1, we assume the training data to have the viewing directions / poses of the single-view semantic-image pairs. However, different from human face, it is difficult to get the poses for cat faces. Considering that a pretrained π -GAN can generate photorealistic cat faces with given random poses, we choose to generate pseudo data for training our Sem2NeRF. Specifically, we randomly sample 30,000 vectors $\mathbf{z} \in \mathbb{R}^{256}$ from the input distribution of π -GAN to generate 30,000 corresponding cat faces by using the released pretrained model⁴. Each image comes with one viewing direction, randomly sampled from normal distributions, with $X \sim \mathcal{N}(\pi/2, 0.3^2)$ for the yaw axis, $X \sim \mathcal{N}(\pi/2, 0.1^2)$ for the pitch axis, and 0 for the roll axis. Camera FOV is set to 18 to ensure the generated image covering the full cat face.

⁴ <https://github.com/marcoamonteiro/pi-GAN>



Fig. 10. CatMask dataset. Left: label legends of 6 semantic classes. Right: single-view cat faces rendered by π -GAN (top row) and the corresponding semantic masks labelled by DatasetGAN (bottom row). Best view in high quality color image

Ray sampling resolution is set to 512×512 , with ray depth range $[0.8, 1.2]$ and ray step size 72. Hierarchical sampling is enabled to improve image quality. We save the rendering viewing direction and the generated images for the CatMask dataset.

Cat semantic by DatasetGAN. A suitable training dataset for Sem2NeRF should be able to be modelled by existing NeRF-based generator, while also comes with semantic labels. However, most datasets used by existing 3D-aware generative models, *e.g.*, cat and car, do not contain component-level semantic masks. We further find out that DatasetGAN can create reasonable semantic masks labels for our task.

DatasetGAN is introduced to automatically build datasets of high-quality semantically segmented images. Specifically, it proposes a MLP-based “Style Interpreter” that can be trained to decode the feature maps of a pretrained StyleGAN model to semantic labels, requiring only very few manually-labeled training samples, *e.g.*, 30 labelled images for cat dataset. In this case, dataset can be automatically built by first randomly sampling images from StyleGAN, following by parsing with the trained style interpreter to obtain corresponding semantic labels.

We use the released⁵ pretrained style interpreter for cat to generate a dataset of 10,000 images-semantic pairs. Such a dataset is further leveraged to train a Deeplab-V3 [5] model, which takes a cat image as input and outputs the corresponding cat semantic mask. We then use the trained Deeplab-V3 to label our generated CatMask dataset. 6 classes are selected based on the label quality. Label legends and examples of the CatMask dataset are given in Figure 10.

A.3 Additional Implementation Details

Style codes averages $(\bar{\gamma}, \bar{\beta})$. We randomly sample 10,000 vectors $\mathbf{z} \in \mathbb{R}^{256}$ from a standard normal distribution, then feed them through the pretrained mapping network of the original π -GAN model, finally average the outputs over the batch dimension to obtain $\bar{\gamma}, \bar{\beta}$.

Datasets. For CelebAMask-HQ [27], both images and semantic masks are loaded as resolution 640×640 , then center cropped to 512×512 . For CatMask,

⁵ https://github.com/nv-tlabs/datasetGAN_release

images and semantic masks are directly loaded as 512×512 . Semantic masks are transformed using one-hot encoding, augmented with the aforementioned contours and distance field representations. We do not apply any other data augmentation, *e.g.*, random flip, to avoid harming the pose information.

Training. Patch discriminator with input size 128×128 is adopted in our experiments, using the implementation provided by the GRAF [46] project⁶. Images are rendered via only the “coarse” network of the decoder, *i.e.*, removing the hierarchical sampling. The sampling range of the scaling factor α of Eq. (2) is initialized as $[0.9, 1.0]$, where the lower bound is exponentially annealed to 0.06 during training. Encoder is initialized with the ImageNet-1K [11] pretrained weights, decoder is initialized with π -GAN pretrained weights, while the discriminator is randomly initialized. We freeze the parameters of the decoder, and set the learning rate of the encoder and discriminator to 1×10^{-4} and 2×10^{-5} , respectively. Ranger optimizer⁷ is used for both encoder and discriminator. We set the training batch size to 8, and use V100 GPUs to train all related models for 200,000 iterations.

Inference. To render qualitative results, rays are cast with size 512×512 and depth step 72 in the inference. Besides, “fine” network is activated to enable hierarchical sampling. For GAN-inversion used by pix2pixHD as mentioned in Section 4.1, we adopt the implementation from the π -GAN project⁸, and set the iteration number to 700 as suggested.

B Additional Visual Results

In the following three pages, we will present additional visual results for the proposed Sem2NeRF regarding free-viewpoint image generation (see Section B.1), semantic mask editing (see Section B.2) and multi-modal synthesis (see Section B.3). Results are demonstrated on both CelebAMask-HQ and CatMask, and they are all best viewed in high quality color image.

⁶ <https://github.com/autonomousvision/graf>

⁷ <https://github.com/lessw2020/Ranger-Deep-Learning-Optimizer>

⁸ <https://github.com/marcoamonteiro/pi-GAN>

B.1 Free-viewpoint Image Generation

Additional visual results of free-viewpoint image generation on CelebAMask-HQ and CatMask datasets are given in Figure 11 and Figure 12, respectively.

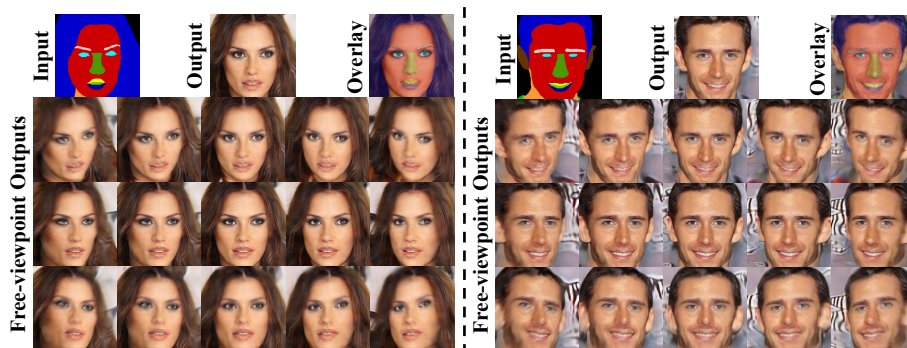


Fig. 11. Free-viewpoint image generation on CelebAMask-HQ. “Output” refers to the generated image that has the same viewing direction as the “Input”, and “Overlay” shows the results of overlaying “Output” with “Input”, so as to better demonstrate the mapping accuracy



Fig. 12. Free-viewpoint image generation on CatMask. “Output” refers to the generated image that has the same viewing direction as the “Input”, and “Overlay” shows the results of overlaying “Output” with “Input”, so as to better demonstrate the mapping accuracy

B.2 Semantic Mask Editing

Additional visual results of semantic mask editing on CelebAMask-HQ and CatMask datasets are given in Figure 13 and Figure 14, respectively.

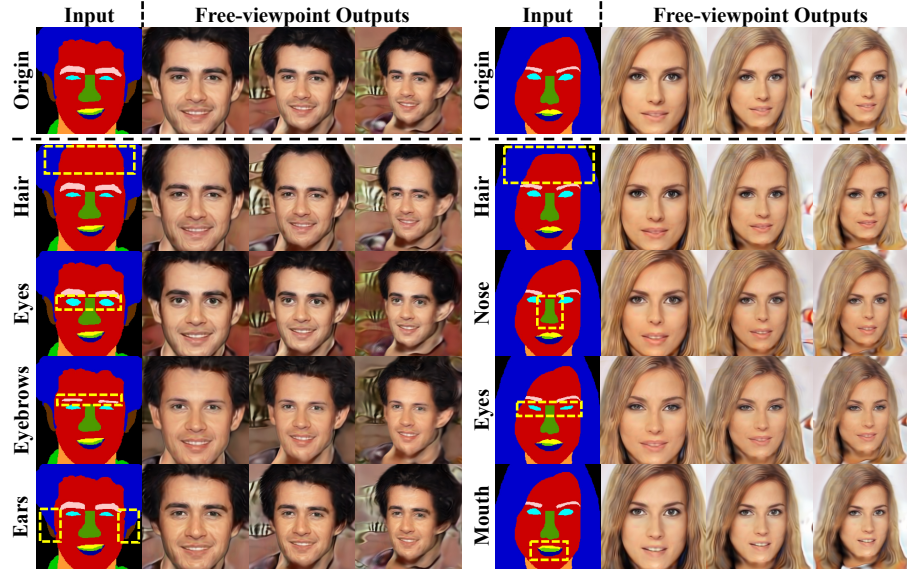


Fig. 13. Semantic mask editing on CelebAMask-HQ. The first row shows the results of the original semantic masks, while the following rows give the results of editing the mentioned area, highlighted with yellow-dash box. Three viewpoints are given for each group, with the first one having the same viewing direction as the input

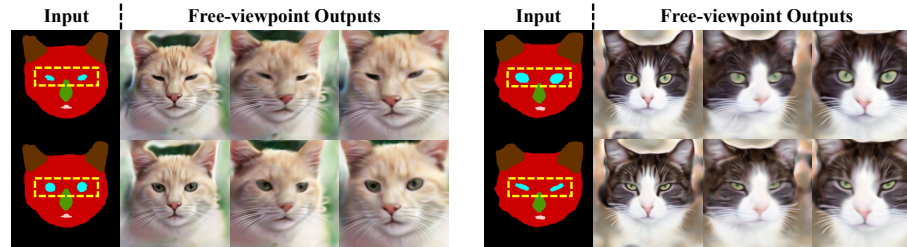


Fig. 14. Semantic mask editing on CatMask. Edited regions are highlighted with yellow-dash box. The first viewpoint has the same pose as the input

B.3 Multi-modal Synthesis

Additional visual results regarding multi-modal synthesis on CelebAMask-HQ and CatMask datasets are given in Figure 15 and Figure 16, respectively.

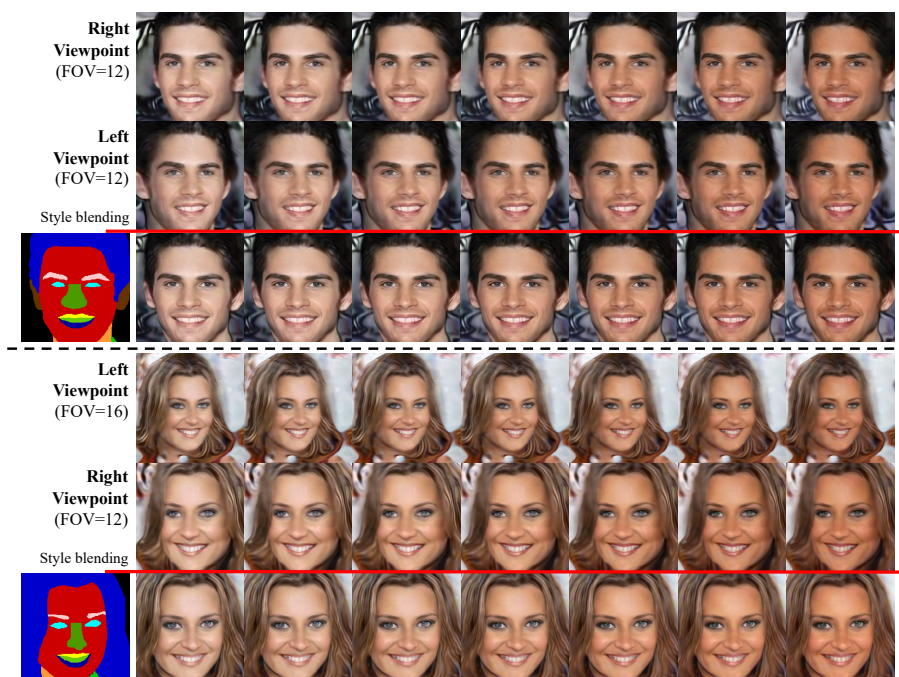


Fig. 15. Multi-modal synthesis on CelebAMask-HQ. Styles are linearly blended from left to right. The last viewpoint in each group has the same pose as the input

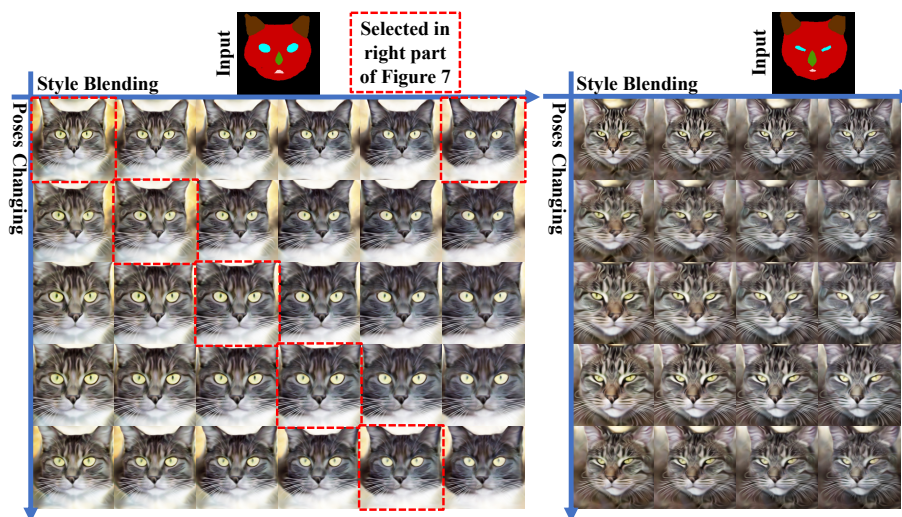


Fig. 16. Multi-modal synthesis on CatMask. Left case shows the full version of Figure 7 (right part), where the selected images are highlighted in red-dash border. Images in the first row have the same viewing direction as the input