

SINE: Semantic-driven Image-based NeRF Editing with Prior-guided Editing Field

Chong Bao^{1*} Yinda Zhang^{2*} Bangbang Yang^{1*} Tianxing Fan¹
Zesong Yang¹ Hujun Bao¹ Guofeng Zhang^{1†} Zhaopeng Cui^{1†}

¹State Key Lab of CAD&CG, Zhejiang University ²Google

<https://zju3dv.github.io/sine/>



Figure 1. We propose a novel semantic-driven image-based editing approach, which allows users to edit a photo-realistic NeRF with a single-view image or with text prompts, and renders edited novel views with vivid details and multi-view consistency.

Abstract

Despite the great success in 2D editing using user-friendly tools, such as Photoshop, semantic strokes, or even text prompts, similar capabilities in 3D areas are still limited, either relying on 3D modeling skills or allowing editing within only a few categories. In this paper, we present a novel semantic-driven NeRF editing approach, which enables users to edit a neural radiance field with a single image, and faithfully delivers edited novel views with high fidelity and multi-view consistency. To achieve this goal, we propose a prior-guided editing field to encode fine-grained geometric and texture editing in 3D space, and develop a series of techniques to aid the editing process, including cyclic constraints with a proxy mesh to facilitate geometric supervision, a color compositing mechanism to stabilize semantic-driven texture editing, and a feature-cluster-based regularization to preserve the irrelevant content unchanged. Extensive experiments and editing examples on both real-world and synthetic data demonstrate that our method achieves photo-realistic 3D editing using only a single edited image, pushing the bound of semantic-driven editing in 3D real-world scenes.

1. Introduction

Semantic-driven editing approaches, such as stroke-based scene editing [36, 41, 70], text-driven image synthesis and editing [1, 53, 56], and attribute-based face editing [28, 64], have greatly improved the ease of artistic creation. However, despite the great success of 2D image edit-

ing and neural rendering techniques [14, 44], similar editing abilities in the 3D area are still limited: (1) they require laborious annotation such as image masks [28, 75] and mesh vertices [73, 78] to achieve the desired manipulation; (2) they conduct global style transfer [12, 13, 16, 21, 79] while ignoring the semantic meaning of each object part (e.g., windows and tires of a vehicle should be textured differently); (3) they can edit on categories by learning a textured 3D latent representation (e.g., 3D-aware GANs with faces and cars etc.) [6, 8, 9, 18, 48, 60, 63, 64], or at a coarse level [37, 68] with basic color assignment or object-level disentanglement [32], but struggle to conduct texture editing on objects with photo-realistic textures or out-of-distribution characteristics.

Based on this observation, we believe that, on the way toward semantic-driven 3D editing, the following properties should be ensured. First, the operation of editing should be effortless, *i.e.*, users can edit 3D scenes on a single 2D image in convenient ways, *e.g.*, using off-the-shelf tools such as GAN-based editing [29, 36], text-driven editing [1, 56], Photoshop, or even a downloaded Internet image without pixel-wise alignment, rather than steering 3D modeling software with specific knowledge [73], or repeatedly editing from multi-view images. Second, the editing method should be applicable to real-world scenes or objects and preserve vivid appearances, which is beyond existing 3D-aware generative models [8, 9] due to the limited categories and insufficient data diversity on real-world objects.

To fulfill this goal, we propose a novel Semantic-driven Image-based Editing approach for Neural radiance field in real-world scenes, named SINE. Specifically, our method allows users to edit a neural radiance field with a sin-

* Authors contributed equally.

† Corresponding authors.

gle image, *i.e.*, either by changing a rendered image using off-the-shelf image editing tools or providing an image for texture transferring (see Sec. 4.4), and then delivers edited novel views with consistent semantic meaning. Unlike previous works that directly fine-tune the existing NeRF model [32, 37, 68], SINE learns a prior-guided editing field to encode geometric and texture changes over the original 3D scene (see Fig. 2), thus enabling fine-grained editing ability. By leveraging guidance from existing neural priors (shape prior models [15] and Vision Transformer models [7], *etc.*), SINE can directly perform semantic-driven editing on photo-realistic scenes without pre-training a category-level latent space. For example, in Fig. 1, users can stretch a car’s back or change all four tires to cookies by only editing a single image, and can even cooperate with text-prompts editing [1] to modify a specific object of a scene with vivid appearances.

However, even when guided with neural priors, editing NeRF from a single image with multi-view consistency and accuracy is still challenging. (1) The generic NeRF does not necessarily provide an explicit surface or signed distance field, such that it cannot directly work with shape priors [15]. Therefore, we propose to use cyclic constraints with a proxy mesh to represent the edited NeRF’s geometry, which facilitates guided editing using coarse shape prior. (2) Learning a coordinate-based 3D editing field using a single edited view is not sufficient to capture fine-grained details, and applying semantic supervision [7, 55] directly to the editing field leads to sub-optimal convergence (see Sec. 4.5). To tackle these challenges, we propose a color compositing mechanism by first rendering the template NeRF color and modification color individually, and then deferred blending them to yield the edited view, which significantly improves semantic-driven texture editing. (3) Ideally, a user’s editing should only affect the desired regions while maintaining other parts untouched. However, in semantic-driven editing, the prior losses require taking the full shape or image as input, which leads to appearance or shape drifting at the undesired area. To precisely control the editing while excluding irrelevant parts from being affected, we generate feature clusters of the editing area using the ViT-based feature field [7, 32], and use these clusters to distinguish whether a location is allowed to be edited or should remain unchanged.

In summary, the contributions of our paper are as follows. (1) We propose a novel semantic-driven image-based NeRF editing approach, called SINE, which allows users to edit a neural radiance field simply on just a single view of the rendering. SINE leverages a prior-guided editing field to encode fine-grained geometry and texture changes over the given pre-trained NeRF, thus delivering multi-view consistent edited views with high fidelity. (2) To achieve semantic editing functionality, we develop a series of techniques,

including cyclic constraints with a proxy mesh for geometric editing, the color compositing mechanism to enhance texture editing, and the feature-cluster-based regularization to control the affected editing area and maintain irrelevant parts unchanged. (3) Experiments and editing examples on both real-world/synthetic and object-centric/unbounded 360° scenes data demonstrate superior editing capabilities and quality with effortless operations.

2. Related Works

Neural rendering with external priors. Neural rendering techniques aim at rendering novel views with high-quality [44] or controllable properties [28, 51] by learning from 2D photo capture. Recently, NeRF [44] achieves photo-realistic rendering with volume rendering and inspires many works, including surface reconstruction [33, 69, 77], scene editing [4, 19, 71, 74, 75] and generation [23, 54], inverse rendering [5, 80], SLAM [76, 81], *etc.* For learning from few-shot images [24] or 3D inpainting [45], NeRF’s variants use hand-crafted losses [47] or large language-image models [24, 72] as external priors. However, due to insufficient 3D supervision, such methods cannot reconstruct accurate geometry and only produce visually plausible results. Besides, some works [22, 34, 42] use the symmetric assumption to reconstruct category-level objects (*e.g.*, cars, chairs) but cannot generalize on complex scenes.

Neural 2D & 3D scene editing. With the development of neural networks, semantic-driven 2D photo editing allows user editing in various friendly ways, such as controlling attribute of faces [20, 29], stroke-based editing [36, 41, 70], sketch-to-image generation [11, 58], image-to-image texture transferring [66], or text-driven image generation [56] and editing [31]. Nevertheless, in 3D scene editing, similar capabilities are still limited due to the high demand for multi-view consistency. Existing approaches either rely on laborious annotation [28, 73, 75, 78], only support object deformation or translation [32, 65, 67, 78], or only perform global style transfer [12, 13, 16, 21, 79] without strong semantic meaning. Recently, 3D-aware GANs [8, 9, 18, 25, 48, 60, 63] and semantic NeRF editing [37, 68] learn a latent space of the category and enable editing via latent code control. However, the quality and editing ability of these methods mainly depend on the dataset (*e.g.*, human faces [63, 64] or objects in ShapeNet [10]), and they cannot generalize to objects with rich appearances or out-of-distribution features [1]. In contrast, our method allows for semantic-driven editing directly on the given photo-realistic NeRF, and uses a prior-guided editing field to learn fine-grained editing from only a single image.

3. Method

We first formulate the goal of our semantic NeRF editing task as follows. As illustrated in the left part of Fig. 2,

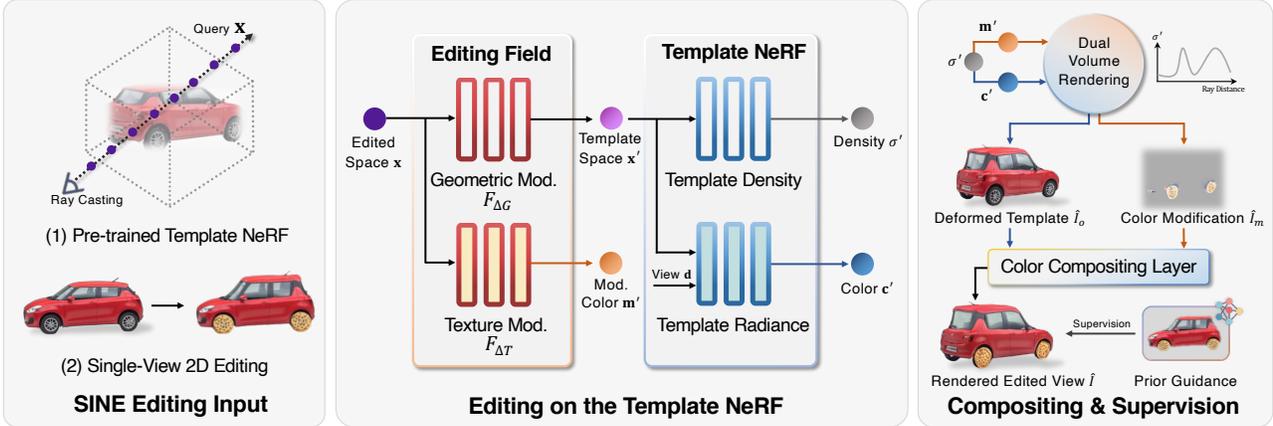


Figure 2. **Overview.** We encode geometric and texture changes over the original template NeRF with a prior-guided editing field, where the geometric modification field $F_{\Delta G}$ transformed the edited space query \mathbf{x} into the template space \mathbf{x}' , and the texture modification field $F_{\Delta T}$ encodes modification colors \mathbf{m}' . Then, we render deformed template image \hat{I}_o and color modification image \hat{I}_m with all the queries, and use a color compositing layer to blend \hat{I}_o and \hat{I}_m into the edited view \hat{I} .

given a pre-trained NeRF of a photo-realistic scene (named template NeRF), we aim at editing the template NeRF using only a single-view 2D image, and then produce novel views with consistent semantic meaning (see Sec. 4.3 and Sec. 4.4). Note that naïvely fine-tuning on the edited single view cannot obtain satisfactory results due to the spatial ambiguity and lack of multi-view supervision (see Sec. 4.2). Therefore, we propose to use a novel prior-guided editing field to encode fine-grained changes (Sec. 3.1) in 3D space, which leverages geometry and texture priors to guide the learning of semantic-driven editing (Sec. 3.2 and Sec. 3.3). Besides, to precisely control the editing area while maintaining other parts unchanged, we design editing regularization with feature cluster-based semantic masking (Sec. 3.4).

3.1. SINE Rendering Pipeline

As illustrated in Fig. 2, we use a dedicated editing field to encode geometry and texture changes over the pre-trained template NeRF. The editing field consists of an implicit geometric modification field $F_{\Delta G}$ and a texture modification field $F_{\Delta T}$, where $F_{\Delta G}$ deforms the query points from the observed edited space to the original template space, as $\mathbf{x}' := F_{\Delta G}(\mathbf{x})$, and $F_{\Delta T}$ encodes the modification color \mathbf{m}' , as $\mathbf{m}' := F_{\Delta T}(\mathbf{x})$. Specifically, for each sampled query point $\{\mathbf{x}_i | i = 1, \dots, N\}$ along the ray \mathbf{r} with view direction \mathbf{d} , we first obtain the deformed points \mathbf{x}' (in template space) and modification color \mathbf{m}' , and feed \mathbf{x}' and \mathbf{d} to the template NeRF to obtain the density δ' and template colors \mathbf{c}' . Then, we perform dual volume rendering both on edited fields and template NeRF following the quadrature rules [40, 44], which is defined as:

$$\begin{aligned} \hat{C}_o(\mathbf{r}) &= \sum_{i=1}^N T_i \alpha_i \mathbf{c}'_i, & \hat{C}_e(\mathbf{r}) &= \sum_{i=1}^N T_i \alpha_i \mathbf{m}'_i, \\ T_i &= \exp\left(-\sum_{j=1}^{i-1} \sigma'_j \delta_j\right), \end{aligned} \quad (1)$$

where $\alpha_i = 1 - \exp(-\sigma'_i \delta_i)$, and δ_i is the distance between adjacent samples along the ray. In this way, we obtain the deformed template image \hat{I}_o from the template NeRF's pixel color $\hat{C}_o(\mathbf{r})$ and color modification image \hat{I}_m from the modification color $\hat{C}_e(\mathbf{r})$. Finally, we apply the color compositing layer (see Sec. 3.3) to blend \hat{I}_o and \hat{I}_m into the resulting edited views \hat{I} .

3.2. Prior-Guided Geometric Editing

In this section, we explain how to learn $F_{\Delta G}(\mathbf{x})$ with the geometric prior.

Shape prior constraint on the edited NeRF. We leverage geometric prior models, such as neural implicit shape representation [15, 49] or depth prediction [3], to mitigate the ambiguity of geometric editing based on editing from a single perspective. (1) For objects within a certain shape category (e.g., cars, airplanes), we use DIF [15], in which the implicit SDF field and the prior mesh \hat{M}_P can be generated with the condition of an optimizable latent code $\hat{\mathbf{z}}$. We force the edited NeRF's geometry \hat{M}_E to be explainable by a pre-trained DIF model with the geometric prior loss:

$$\begin{aligned} \mathcal{L}_{gp} &= \min_{\hat{\mathbf{z}}} \left(\sum_{\mathbf{p}' \in \hat{M}_E} f_{\text{SDF}}(\hat{\mathbf{z}}, \mathbf{p}') + \lambda \|\hat{\mathbf{z}}\|_2^2 \right) \\ &+ \sum_{\mathbf{p}'_i \in \hat{M}_E} \min_{\mathbf{p}_t \in \hat{M}_P} \|\mathbf{p}'_i - \mathbf{p}_t\|_2^2 \\ &+ \sum_{\mathbf{p}_i \in \hat{M}_P} \min_{\mathbf{p}'_t \in \hat{M}_E} \|\mathbf{p}_i - \mathbf{p}'_t\|_2^2. \end{aligned} \quad (2)$$

The first term encourages the sampled surface points on the edited NeRF's geometry \hat{M}_E to lie on the manifold of DIF's latent space with an SDF loss f_{SDF} and the latent code regularization [15]. The last two terms are Chamfer constraints, which enforce the \hat{M}_E close to the DIF's periodically updated prior mesh \hat{M}_P [38] by minimizing the closest surface points. (2) For objects without a category-level prior, we can build a finalized shape prior \hat{M}_P beforehand. Practically, we find 3D deforming vertices with 2D correspon-

dence [27] and monocular depth prediction [3], and use ARAP [61] to deform the proxy triangle mesh M_Θ to \hat{M}_P . Then, we can inherit the Chamfer loss term in Eq. (2) for prior-guided supervision.

Representing edited NeRF’s geometry as a deformed proxy mesh. The edited NeRF has no explicit surface definition or SDF field to directly apply the geometric prior loss (Eq. (2)). Therefore, to obtain the edited mesh surface \hat{M}_E , as illustrated in Fig. 3 (a), we first fit the template NeRF geometry with a proxy mesh M_Θ [38, 69], and then learn a forward modification field $F'_{\Delta G}$ to warp the template proxy mesh to the edited space. $F'_{\Delta G}$ is an inverse of the editing field $F_{\Delta G}$, which maps from the template space to the query space [35, 43], and can be supervised using a cycle loss \mathcal{L}_{cyc} (see the supplementary material for details). Note that the deformed mesh proxy might not reflect fine-grained details of the specific shape identity. It facilitates applying shape priors to the edited field and provides essential guidance during geometric editing.

Learning geometric editing with users’ 2D editing. The goal of geometric editing is to deform the given NeRF according to the edited target image while satisfying semantic properties. To this end, apart from the geometric prior loss in Eq. (2), we add the following geometric editing loss in two folds. **(1)** We encourage the edited NeRF to satisfy the user’s edited image by directly supervising rendering colors and opacity on N_r rays, which is defined as:

$$\mathcal{L}_{gt} = \frac{1}{|N_r|} \sum_{\mathbf{r} \in N_r} \|\hat{C}(\mathbf{r}) - C_t(\mathbf{r})\|_2^2 + \text{BCE}(\hat{O}(\mathbf{r}), O_e(\mathbf{r})). \quad (3)$$

The first photometric loss term encourages the rendered color $\hat{C}(\mathbf{r})$ close to the edited target color $C_t(\mathbf{r})$. The second silhouette loss term enforces the rendered opacity $\hat{O}(\mathbf{r})$ close to the edited object’s silhouette $O_e(\mathbf{r})$ (derived from users’ editing tools) by minimizing the binary cross-entropy loss, where $\hat{O}(\mathbf{r}) = \sum_{i=1}^N T_i \alpha_i$. **(2)** To obtain a spatially smooth deformation and mitigate overfitting to the mesh’s surface points, inspired by previous works [15, 51, 52], we also add deformation regularization as:

$$\mathcal{L}_{gr} = \frac{1}{M} \sum_{i=1}^N \|\nabla F_{\Delta G}(\mathbf{p}_i)\|_2 + \|F_{\Delta G}(\mathbf{p}_i) - F_{\Delta G}(\mathbf{p}_i + \epsilon)\|_1, \quad (4)$$

where the first term penalizes the spatial gradient of the geometric editing, and the second term encourages the editing to be smooth under a mild 3D positional jitter ϵ .

The overall geometric editing loss is defined as:

$$\mathcal{L}_{geo} = \lambda_{gp} \mathcal{L}_{gp} + \mathcal{L}_{gt} + \lambda_{gr} \mathcal{L}_{gr} + \lambda_{cyc} \mathcal{L}_{cyc}, \quad (5)$$

where we set $\lambda_{gp} = 0.03$, $\lambda_{gr} = 0.1$ and $\lambda_{cyc} = 10$. Intuitively, the geometric editing loss \mathcal{L}_{geo} jointly optimizes edited NeRF’s geometry \hat{M}_E and the latent shape code $\hat{\mathbf{z}}$ (for category-level objects) to best fit the user’s 2D editing

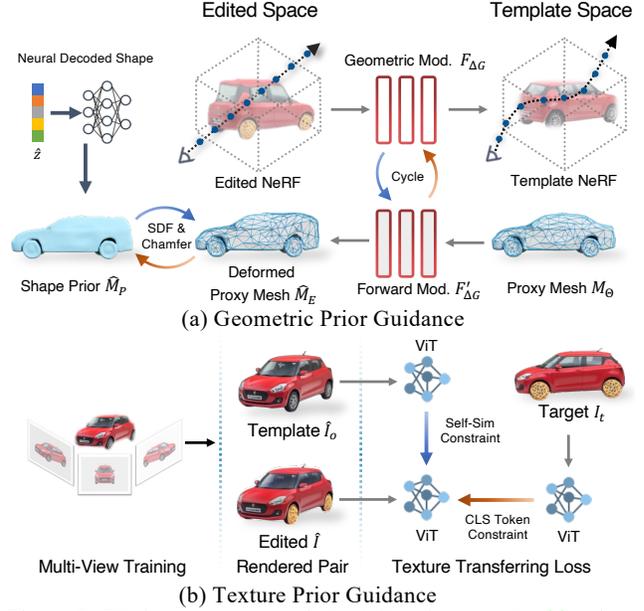


Figure 3. We leverage geometric [3, 15] and texture [7, 66] priors to guide the learning of semantic-driven NeRF editing.

while maintaining shape prior’s semantic properties (*i.e.*, shape symmetry or physical conformity).

3.3. Prior-Guided Texture Editing

Semantic texture prior supervision. In our task, users only conduct editing on a single image, but we hope to naturally propagate editing effects to multi-views with semantic meaning (see Fig. 1). Therefore, we need to utilize semantic texture supervision that supports transferring the editing to the given NeRF across views, rather than using a pixel-aligned photometric loss. Inspired by Tumanyan *et al.* [66], we use a pre-trained ViT model [7] as the semantic texture prior, and apply the texture transferring loss in a multi-view manner as illustrated in Fig. 3 (b), which is defined as:

$$\mathcal{L}_{tex} = \|t_{CLS}(I_t) - t_{CLS}(\hat{I})\|_2 + \|S(\hat{I}_o) - S(\hat{I})\|_F, \quad (6)$$

where I_t is the user’s edited image, \hat{I}_o and \hat{I} are the template image and edited image as introduced in Sec. 3.1. $t_{CLS}(\cdot)$ and $S(\cdot)$ are the extracted deepest CLS token and the structural self-similarity defined by Tumanyan *et al.* [66]. Essentially, this loss encourages \hat{I}_o and \hat{I} to share a similar spatial structure, and \hat{I}_t and \hat{I} to contain similar image cues.

Decoupled rendering with color compositing layer. To achieve texture modification, a naïve approach is to directly add the modification color \mathbf{m}' from the editing field to the template NeRF’s radiance color \mathbf{c}' during volume rendering. However, we find it suffers from sub-optimal convergence when cooperating with texture transferring loss (see Sec. 4.5), since NeRF struggles to learn the global-consistent appearance under the variational supervisory as shown in Fig. 8(a). To tackle this issue, we re-design the rendering pipeline in a decoupled manner. As shown in Fig. 2, we first render the deformed template image \hat{I}_o with

template NeRF and the color modification \hat{I}_m with $F_{\Delta T}$, and then use a 2D CNN-based color compositing layer to deferred blend the modification \hat{I}_m into the template image \hat{I}_o , which yields final edited view \hat{I} . Intuitively, the coordinate-based editing field can encode fine-grained details from photometric constraints but cannot easily learn from coarse semantic supervision, while the proposed color compositing layer can reduce the difficulty by using easy-to-learn CNN layers before applying texture transferring loss. Besides, it also learns view-dependent effects from the semantic prior, making the rendering results more realistic (e.g., the shining diamond effect in Fig. 1).

3.4. Editing Regularization

Feature-cluster-based semantic masking. To precisely edit the desired region while preserving other content unchanged, inspired by previous works [32, 65, 67], we learn a distilled feature field with DINO-ViT [7] to reconstruct scenes/objects with semantic features. However, existing semantic field decomposing approaches [32, 65] are limited to the query-based similarity and require all the editing to be finalized on the 3D field, which is not compatible with our color compositing mechanism. Therefore, we leverage users’ editing silhouette M_e to generate several feature clusters from the distilled feature map, and compute semantic masks \hat{M}_e using the closest cosine similarity to cluster centers with a threshold, which will be served for image-based editing regularization.

Regularization on geometric and texture editing. With the semantic masks that indicate the editing area, we can apply editing regularization to the geometric and texture editing, *i.e.*, by enforcing the rendered pixels and the queries at the irrelevant part unchanged, which is defined as:

$$\mathcal{L}_{\text{reg}} = \sum_{\mathbf{x} \in \hat{I} \setminus \hat{M}_e} \|F_{\Delta G}(\mathbf{x})\|_1 + \sum_{\mathbf{r} \in \hat{I} \setminus \hat{M}_e} \|\hat{C}(\mathbf{r}) - \hat{C}_o(\mathbf{r})\|_2^2, \quad (7)$$

where the sampled points \mathbf{x} and rays \mathbf{r} are both from the background area of the computed semantic masks \hat{M}_e .

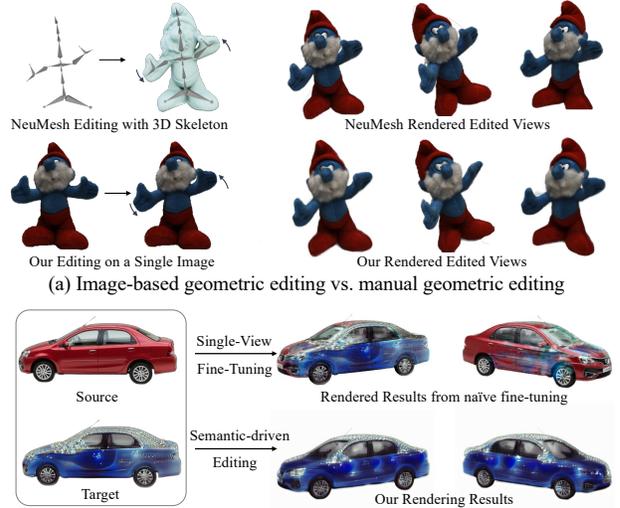
4. Experiments

4.1. Datasets

We evaluate SINE on both real-world/synthetic and object/scene datasets, including real-world car dataset [57], PhotoShape datasets (synthetic chairs) [50], “pinecone”, “vasedeck”, and “garden” from NeRF real-world 360° scenes [2, 44], “chairs” and “hotdog” from NeRF photo-realistic synthetic data [44], bird status from DTU [26] dataset, and “airplane” from BlenderSwap [39]. Please refer to the supplementary material for more details.

4.2. Semantic-driven vs. Manual Editing

We first clarify the difference between our semantic-driven NeRF editing and manual NeRF editing (e.g., NeuMesh [73], NeRF-Editing [78]). As illustrated in Fig. 4

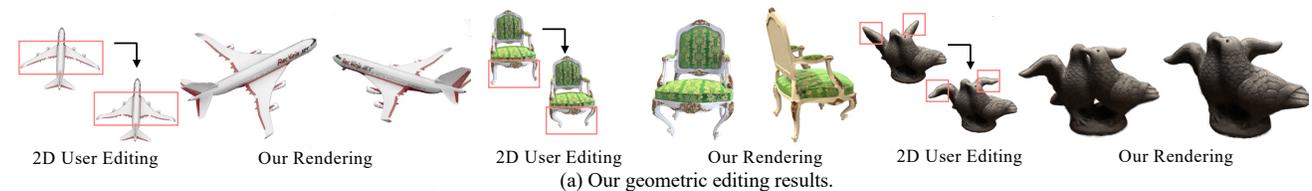


(a) Image-based geometric editing vs. manual geometric editing
 (b) Semantic-driven texture editing vs. naive single-view fine-tuning
 Figure 4. We show the difference between our semantic-driven image-based NeRF editing and manual NeRF editing [73].

(a), our method provides much more effortless ways than manual approaches. For example, they require 3D modeling skills to bind the skeleton of the mesh using Blender [73], or drive models [78] with Mixamo poses [62], while our method can easily achieve similar geometric editing with only a single-view image. Besides, naively fine-tuning NeRF on a single-view with a pixel-aligned photometric loss like NeuMesh [73] would only modify visible regions, which leads to inconsistent novel view rendering (e.g., in Fig. 4 (b), the car edited by single-view fine-tuning would expose unpainted red part). On the contrary, our method leverages semantic priors [7] to naturally edit objects with multi-view consistency, which does not require pixel-wise alignment and enables texture transferring between objects with different shapes (see cars and chairs in Fig. 6 (a)).

4.3. Semantic-driven Geometric Editing

We first show our geometric editing results in Fig. 5 (a), where the objects can be faithfully deformed according to users’ 2D editing (e.g., the airplane with warped wings [39], green chair with bent legs [44] and deformed bird status [26]). For the usage of geometry prior, we use a pre-trained DIF [15] model for cars [57], chairs [50] and planes [39], and use ARAP-based shape priors for general objects without a category-level prior (*i.e.*, toy in Fig. 4 (a), green chair with unusual shape and status in Fig. 5 (a), *etc.*). Then, we compare our method with EG3D [9], a 3D-aware generative model that learns a latent representation of category-level objects, and EditNeRF [37], a NeRF-variant that supports object editing with single-view user interaction. In addition to editing comparisons, we also used PSNR, SSIM, and LPIPS [44] to measure the edited rendering quality on synthetic cars and chairs. For the geometric editing on EG3D, we first conduct 3D GAN inversion to obtain the style code with multi-view images (same input

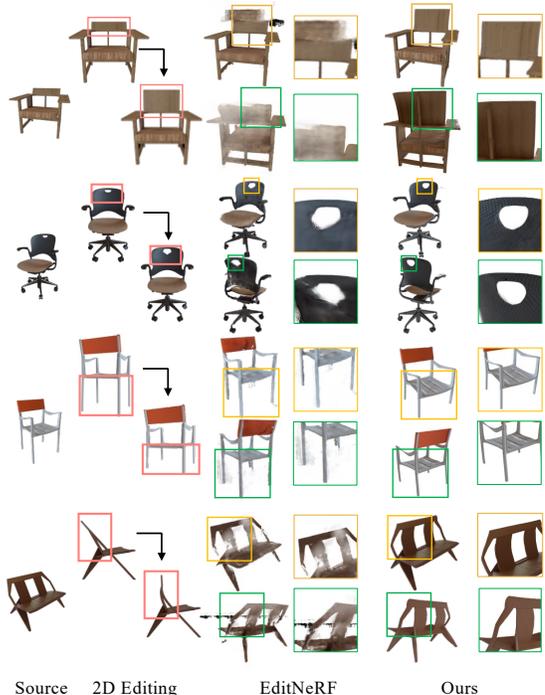


(a) Our geometric editing results.



Source 2D Editing EG3D Ours
 PNSR \uparrow / SSIM \uparrow / LPIPS \downarrow 18.66 / 0.87 / 0.10 33.99 / 0.98 / 0.01

(b) Comparison with EG3D on the real-world car dataset.



Source 2D Editing EditNeRF Ours
 21.23 / 0.92 / 0.16 24.96 / 0.94 / 0.06

(c) Comparison with EditNeRF on the PhotoShape dataset.

Figure 5. We compare the geometric editing with EG3D [9] and EditNeRF [37] on the real-world cars [57] and PhotoShape [50].

as ours), and then fine-tune the code on the target images. As shown in Fig. 5 (b), we conduct different editing operations on four cars from CalWare 360° datasets with DIF shape prior [15], *i.e.*, enlarging/shrinking tires/back. Due to the difficulty of learning a latent textured 3D representation and the limitation of data diversity, 3D-aware generative models like EG3D cannot produce rendering results with fine-grained details, which also results in lower evaluation metrics. For the Photoshape [50], EditNeRF [37] does not provide edited GT images, so we regenerate all testing cases using Blender, which is more challenging than the original ones. Then, we evaluate EditNeRF [37] by fine-tuning the pre-trained models on specific chairs from the PhotoShape dataset. As shown in Fig. 5 (c), EditNeRF produces more blurry rendering results than ours, and cannot achieve satisfactory results with single-view editing (*e.g.*, multi-view inconsistent chair back in the first row, and unmodified or blurry shapes in the third and fourth rows). By contrast, our method consistently delivers high-fidelity rendering results and achieves reliable editing capability by leveraging geometric priors [3, 15]. This demonstrates that, for seman-

tic geometric NeRF editing, learning a prior-guided editing field like ours can maintain better visual quality and achieve greater generalization ability than pre-training a textured 3D generative model or latent model.

4.4. Semantic-driven Texture Editing

We evaluate our semantic texture editing ability on both objects (cars from CalWare 360°, chairs from PhotoShape [50]) and unbounded 360° scenes [2, 44]. Since our method only requires a single image as editing input, we exhibit several editing functionalities as shown in Fig. 6. Users can edit by assigning new textures on the car using Photoshop (adding sea wave windows in Fig. 6 (a)), using a downloaded Internet image with different shapes as a reference (transferring textures of cars and chairs in Fig. 6 (a)). Moreover, we cooperate SINE with off-the-shelf text-prompts editing methods [1] by using a single text-edited image as the target, which enables to change the object’s appearance in the 360° scene with vivid effects (*e.g.*, shiny plastic round table or burning pinecone in Fig.6 (b)) while preserving background unchanged. It is noteworthy that our method does not pre-train a latent model within a spe-

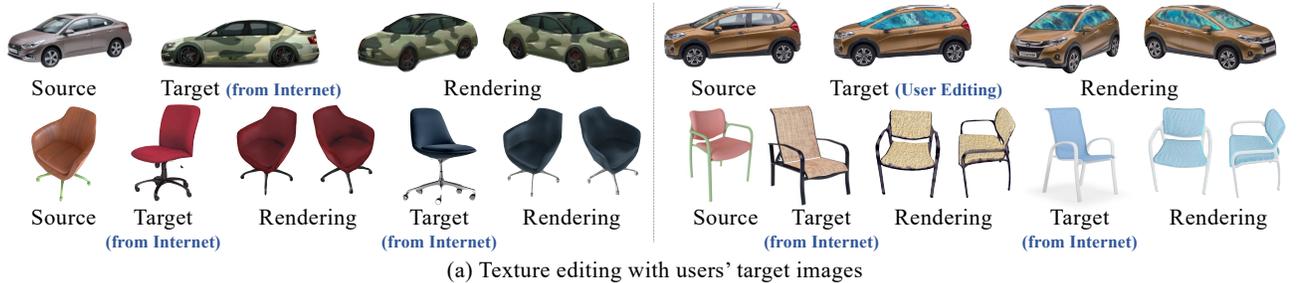


Figure 6. We show our texture editing results when given users' target images and cooperating with text-prompt-based editing methods [1].

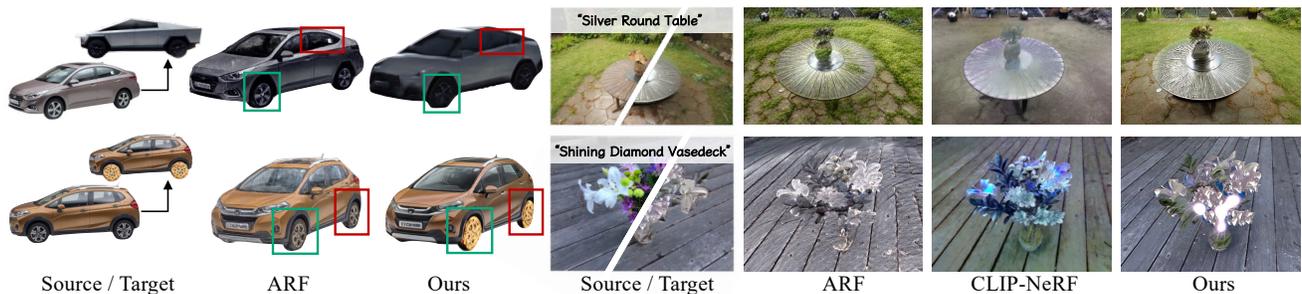


Figure 7. We compare our texture editing with ARF [79] and CLIP-NeRF [68] on the real-world cars [57] and 360° [2, 44] scene dataset.

cific category like cars or chairs, yet still transfers texture between objects with correct semantic meaning, *e.g.*, the texture styles of chair legs and cloths in the edited views are precisely matched to the target images in Fig. 6 (a). Besides, we also compare our methods with ARF [79], a

NeRF stylization method that also takes a single reference image as input, and CLIP-NeRF [68], which supports text-driven NeRF editing using the large language model [55]. As demonstrated in Fig. 7, ARF globally changes appearance colors to the given target images but fails to produce

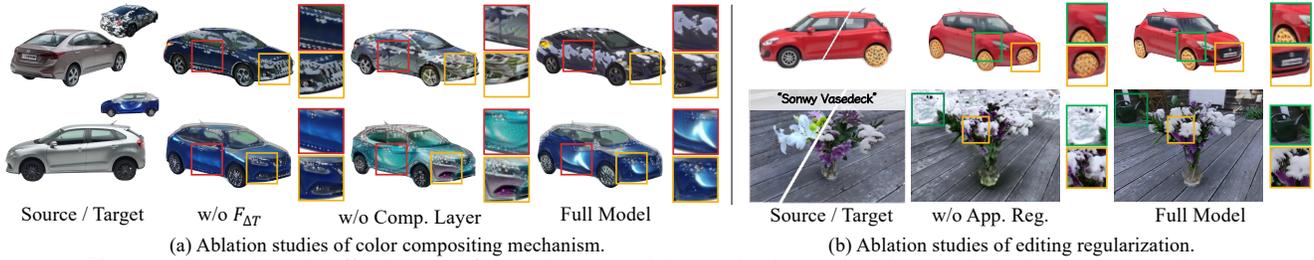


Figure 8. We analyze the effectiveness of the color compositing mechanism and editing regularization in texture editing.

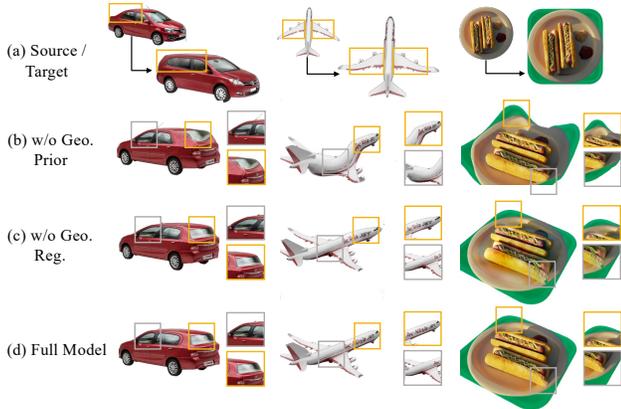


Figure 9. We inspect the efficacy of the geometric prior constraint and editing regularization in geometric editing.

fine-grained details (e.g., cookie tires in Fig. 7 (a)). For CLIP-NeRF, since it directly fine-tunes NeRF’s color layers, the results only show color/hue adjustment on the original scene (e.g., in Fig. 7, the round table turns gray instead of a realistic silver texture, the vase deck turns blue instead of a shining diamond). Thanks to the prior-guided editing field, our method learns more fine-grained editing details than the others, and achieves texture editing with consistent semantic meaning to the given target images (e.g., similar appearance to the Tesla’s cybertruck in Fig. 7 (a)), and delivers rich appearance details and vivid effects (e.g., silver texture and shining diamond effects in Fig. 7 (b)).

User study. We also perform user studies to compare our texture editing (including target-image-based editing and text-prompts editing as Fig. 7) with ARF [79] and CLIP-NeRF [68] on 43 cases with 30 users. The results show that users prefer our methods (89.5% / 83.3%) to ARF (10.4% / 7.4%) and CLIP-NeRF (9.3%). Please refer to the supplementary material for more details.

4.5. Ablation Studies

Geometric prior constraints. We first analyze the effectiveness of geometric prior constraints (Sec. 3.2) by ablating geometric prior loss (Eq. (2)) in Fig. 9 (i.e., deforming cars and airplanes with DIF shape prior, and adding plates with general shape prior). As shown in Fig. 9 (b), when learning without geometric prior constraints, the object will be distorted when rendered from other views (e.g., collapsed car back, twisted airplanes, and warped hotdog plates). By ap-

plying geometric prior constraints, we successfully mitigate the geometric ambiguity for single-image-based editing and produce plausible rendering results from novel views.

Color compositing mechanism. We then inspect the efficacy of the color compositing mechanism (Sec. 3.3) by disabling the texture modification field $F_{\Delta T}$ and the color compositing layer in turn. As demonstrated in Fig. 8 (a), when learning texture editing without $F_{\Delta T}$, the rendered edited object can show a similar global appearance to the target, but lose vivid local patterns (e.g., gray and white grains and blue shininess). When ablating the color compositing layers, the editing effect might not be properly applied to every part of the object (e.g., the uncovered gray part of the car’s front). When all the compositing mechanism is enabled, we successfully learn NeRF editing with fine-grained local patterns and globally similar appearance.

Editing regularization. We finally evaluate the editing regularization (Sec. 3.4) in geometric and texture editing by ablating regularization loss (Eq. (7)). As shown in Fig. 9 (c) and Fig. 8 (b), when learning editing without regularization, the irrelevant part would be inevitably changed (e.g., bent car’s front and airplane’s head in Fig. 9 (c), a spurious cookie at the car’s front and snowy background in Fig. 8 (b)). By adding editing regularization, we can modify the user-desired objects precisely while preserving other content unchanged.

Please refer to the supplementary material for more experiments (e.g., ablation on more loss terms, visualization of color composition layer, discussion with external supervision, etc.).

5. Conclusion

We have proposed a novel semantic-driven NeRF editing approach, which supports editing a photo-realistic template NeRF with a single user-edited image, and deliver edited novel views with high-fidelity and multi-view consistency. As limitation, our approach does not support editing with topology changes, which can be future work. Besides, our method assumes users’ editing to be semantically meaningful, so we cannot use target images with meaningless random paintings.

Acknowledgment. This work was partially supported by NSF of China (No. 62102356).

References

- [1] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. Text2live: Text-driven layered image and video editing. *arXiv preprint arXiv:2204.02491*, 2022. 1, 2, 6, 7, 13, 16, 17
- [2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 5, 6, 7
- [3] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4009–4018, 2021. 3, 4, 6
- [4] Bomng Zhao and Bangbang Yang, Zhenyang Li, Zuoyue Li, Guofeng Zhang, Jiashu Zhao, Dawei Yin, Zhaopeng Cui, and Hujun Bao. Factorized and controllable neural re-rendering of outdoor scene for photo extrapolation. In *Proceedings of the 30th ACM International Conference on Multimedia*, 2022. 2
- [5] Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan Barron, and Hendrik Lensch. Neural-pil: Neural pre-integrated lighting for reflectance decomposition. *Advances in Neural Information Processing Systems*, 34:10691–10704, 2021. 2
- [6] Shengqu Cai, Anton Obukhov, Dengxin Dai, and Luc Van Gool. Pix2nerf: Unsupervised conditional p-gan for single image to neural radiance fields translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3981–3990, June 2022. 1
- [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. 2, 4, 5
- [8] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, 2021. 1, 2
- [9] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 5, 6, 14, 15
- [10] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 2
- [11] Wengling Chen and James Hays. Sketchygan: Towards diverse and realistic sketch to image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9416–9425, 2018. 2
- [12] Yaosen Chen, Qi Yuan, Zhiqiang Li, Yuegen Liu, Wei Wang, Chaoping Xie, Xuming Wen, and Qien Yu. UPST-NeRF: Universal photorealistic style transfer of neural radiance fields for 3d scene. In *arXiv preprint arXiv:2208.07059*, 2022. 1, 2
- [13] Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Wei-Sheng Lai, and Wei-Chen Chiu. Stylizing 3d scene via implicit representation and hypernetwork. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1475–1484, 2022. 1, 2
- [14] Frank Dellaert and Lin Yen-Chen. Neural volume rendering: Nerf and beyond. *arXiv preprint arXiv:2101.05204*, 2020. 1
- [15] Yu Deng, Jiaolong Yang, and Xin Tong. Deformed implicit field: Modeling 3d shapes with learned dense correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10286–10296, 2021. 2, 3, 4, 5, 6, 14
- [16] Zhiwen Fan, Yifan Jiang, Peihao Wang, Xinyu Gong, Dejjia Xu, and Zhangyang Wang. Unified implicit neural stylization. *arXiv preprint arXiv:2204.01943*, 2022. 1, 2
- [17] Kaleido AI GmbH. removebg.bg. <https://www.remove.bg/>, 2022. Accessed: 2022-11-10. 13
- [18] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021. 1, 2
- [19] Michelle Guo, Alireza Fathi, Jiajun Wu, and Thomas Funkhouser. Object-centric neural scene rendering. *arXiv preprint arXiv:2012.08503*, 2020. 2
- [20] Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. Attgan: Facial attribute editing by only changing what you want. *IEEE transactions on image processing*, 28(11):5464–5478, 2019. 2
- [21] Yi-Hua Huang, Yue He, Yu-Jie Yuan, Yu-Kun Lai, and Lin Gao. Stylizednerf: consistent 3d scene stylization as stylized nerf via 2d-3d mutual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18342–18352, 2022. 1, 2
- [22] Eldar Insafutdinov, Dylan Campbell, João F Henriques, and Andrea Vedaldi. Snes: Learning probably symmetric neural surfaces from incomplete data. *arXiv preprint arXiv:2206.06340*, 2022. 2
- [23] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 867–876, 2022. 2
- [24] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. 2
- [25] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12949–12958, 2021. 2
- [26] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large Scale Multi-view Stereopsis Eval-

- uation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 406–413. IEEE, 2014. 5, 15
- [27] Wei Jiang, Eduard Trulls, Jan Hosang, Andrea Tagliasacchi, and Kwang Moo Yi. Cotr: Correspondence transformer for matching across images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6207–6217, 2021. 4
- [28] Kacper Kania, Kwang Moo Yi, Marek Kowalski, Tomasz Trzciński, and Andrea Tagliasacchi. Conerf: Controllable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18623–18632, 2022. 1, 2
- [29] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1, 2
- [30] Yoni Kasten, Dolev Ofri, Oliver Wang, and Tali Dekel. Layered neural atlases for consistent video editing. *ACM Transactions on Graphics (TOG)*, 40(6):1–12, 2021. 16
- [31] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. 2
- [32] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. *arXiv preprint arXiv:2205.15585*, 2022. 1, 2, 5, 15, 17
- [33] Hai Li, Xingrui Yang, Hongjia Zhai, Yuqian Liu, Hujun Bao, and Guofeng Zhang. Vox-surf: Voxel-based implicit surface representation. *IEEE Transactions on Visualization and Computer Graphics*, 2022. 2
- [34] Xingyi Li, Chaoyi Hong, Yiran Wang, Zhiguo Cao, Ke Xian, and Guosheng Lin. Symmnerf: Learning to explore symmetry prior for single-view view synthesis. *arXiv preprint arXiv:2209.14819*, 2022. 2
- [35] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6498–6508, 2021. 4, 14
- [36] Huan Ling, Karsten Kreis, Daiqing Li, Seung Wook Kim, Antonio Torralba, and Sanja Fidler. Editgan: High-precision semantic image editing. *Advances in Neural Information Processing Systems*, 34:16331–16345, 2021. 1, 2, 13, 15
- [37] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5773–5783, 2021. 1, 2, 5, 6, 13
- [38] William E Lorensen and Harvey E Cline. Marching Cubes: A High Resolution 3D Surface Construction Algorithm. *ACM SIGGRAPH Computer Graphics*, 21(4):163–169, 1987. 3, 4
- [39] John Roper Matthew Muldoon. Blenderswap. <https://www.blenderswap.com/>, 2022. Accessed: 2022-11-10. 5, 13, 14
- [40] Nelson L. Max. Optical Models for Direct Volume Rendering. *IEEE Trans. Vis. Comput. Graph.*, 1(2):99–108, 1995. 3
- [41] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 1, 2
- [42] Lu Mi, Abhijit Kundu, David Ross, Frank Dellaert, Noah Snavely, and Alireza Fathi. im2nerf: Image to neural radiance field in the wild. *arXiv preprint arXiv:2209.04061*, 2022. 2
- [43] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. Leap: Learning articulated occupancy of people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10461–10471, 2021. 4, 14
- [44] 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. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2, 3, 5, 6, 7, 13, 16, 17
- [45] Ashkan Mirzaei, Yash Kant, Jonathan Kelly, and Igor Gilitschenski. Laterf: Label and text driven object radiance fields. *arXiv preprint arXiv:2207.01583*, 2022. 2
- [46] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 13
- [47] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. 2
- [48] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 1, 2
- [49] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 3, 14
- [50] Keunhong Park, Konstantinos Rematas, Ali Farhadi, and Steven M. Seitz. Photoshape: Photorealistic materials for large-scale shape collections. *ACM Trans. Graph.*, 37(6), Nov. 2018. 5, 6, 13
- [51] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 2, 4, 13, 14
- [52] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural

- radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. [4](#), [13](#), [18](#)
- [53] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. [1](#)
- [54] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. [2](#)
- [55] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. [2](#), [7](#)
- [56] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. [1](#), [2](#)
- [57] Arun Kumar Sahlam. Carwale. <https://www.carwale.com/>, 2022. Accessed: 2022-11-10. [5](#), [6](#), [7](#), [13](#), [15](#), [17](#)
- [58] Patsorn Sangkloy, Jingwan Lu, Chen Fang, Fisher Yu, and James Hays. Scribbler: Controlling deep image synthesis with sketch and color. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5400–5409, 2017. [2](#)
- [59] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. [13](#)
- [60] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020. [1](#), [2](#)
- [61] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry processing*, volume 4, pages 109–116, 2007. [4](#)
- [62] Nazim Kareemi Stefano Corazza. Mixamo. <https://www.mixamo.com/>, 2022. Accessed: 2022-11-10. [5](#)
- [63] 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. *arXiv preprint arXiv:2205.15517*, 2022. [1](#), [2](#)
- [64] Jingxiang Sun, Xuan Wang, Yong Zhang, Xiaoyu Li, Qi Zhang, Yebin Liu, and Jue Wang. Fenerf: Face editing in neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7672–7682, 2022. [1](#), [2](#)
- [65] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. *arXiv preprint arXiv:2209.03494*, 2022. [2](#), [5](#)
- [66] Narek Tumanyan, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Splicing vit features for semantic appearance transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10748–10757, 2022. [2](#), [4](#), [15](#), [16](#)
- [67] Suhani Vora, Noha Radwan, Klaus Greff, Henning Meyer, Kyle Genova, Mehdi SM Sajjadi, Etienne Pot, Andrea Tagliasacchi, and Daniel Duckworth. Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. *arXiv preprint arXiv:2111.13260*, 2021. [2](#), [5](#)
- [68] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022. [1](#), [2](#), [7](#), [8](#), [15](#), [17](#)
- [69] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. [2](#), [4](#), [14](#)
- [70] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. [1](#), [2](#)
- [71] Qianyi Wu, Xian Liu, Yuedong Chen, Kejie Li, Chuanxia Zheng, Jianfei Cai, and Jianmin Zheng. Object-compositional neural implicit surfaces. *arXiv preprint arXiv:2207.09686*, 2022. [2](#)
- [72] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. *arXiv preprint arXiv:2204.00928*, 2022. [2](#)
- [73] Bangbang Yang, Chong Bao, Junyi Zeng, Hujun Bao, Yinda Zhang, Zhaopeng Cui, and Guofeng Zhang. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision*, pages 597–614. Springer, 2022. [1](#), [2](#), [5](#)
- [74] Bangbang Yang, Yinda Zhang, Yijin Li, Zhaopeng Cui, Sean Fanello, Hujun Bao, and Guofeng Zhang. Neural rendering in a room: amodal 3d understanding and free-viewpoint rendering for the closed scene composed of pre-captured objects. *ACM Transactions on Graphics (TOG)*, 41(4):1–10, 2022. [2](#)
- [75] Bangbang Yang, Yinda Zhang, Yinghao Xu, Yijin Li, Han Zhou, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Learning object-compositional neural radiance field for editable scene rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13779–13788, 2021. [1](#), [2](#)
- [76] Xingrui Yang, Hai Li, Hongjia Zhai, Yuhang Ming, Yuqian Liu, and Guofeng Zhang. Vox-fusion: Dense tracking and mapping with voxel-based neural implicit representation. In *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 499–507. IEEE, 2022. [2](#)
- [77] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *arXiv preprint arXiv:2106.12052*, 2021. [2](#)
- [78] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. Nerf-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18353–18364, 2022. [1](#), [2](#), [5](#)

- [79] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields. In *European Conference on Computer Vision*, pages 717–733. Springer, 2022. [1](#), [2](#), [7](#), [8](#), [14](#), [17](#)
- [80] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021. [2](#)
- [81] 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*, pages 12786–12796, 2022. [2](#)

Supplementary Material

In this supplementary material, we describe more details of our method, including model architecture in Sec. A, dataset preparation in Sec. B, and implementation details in Sec. C. Besides, we also conduct more experiments in Sec. D. More qualitative results can be found in our supplementary video, and the source code will be released upon the acceptance of this paper.

A. Model Architecture

We first explain the details of the model architecture. Specifically, we adopt the multi-resolution voxel-hashing encoder by Müller *et al.* [46] as the coordinate-based encoder, and build the template NeRF and the editing field in a decoupled manner. The voxel-hashing encoder is constructed with 16 levels with 2-dimensional features for each level. For the template NeRF, we use the voxel-hashing encoder to encode the queries' coordinates and use spherical harmonics with 4 degrees to encode the ray direction. The density and color heads for model output consist of 1 hidden layer with 128 hidden size and 2 hidden layers with 64 hidden size, respectively. As introduced in Sec. 3.1, the editing field consists of a geometric modification field $F_{\Delta G}$ and a texture modification field $F_{\Delta T}$. The geometric modification field $F_{\Delta G}$ and the corresponding forward modification field $F'_{\Delta G}$ are both constructed with an MLP of 1 hidden layer and 128 hidden size with the ReLU activation, and we adopt the positional encoding [44] (with 4 frequencies) to all input query points. The texture modification field $F_{\Delta T}$ is constructed with a voxel-hashing encoder (same size as the template NeRF), followed by an MLP of 1 hidden layer and 128 hidden size with ReLU activation. During the dual volume rendering stage, we follow Mildenhall *et al.* [44] by using 64 coarse samples and 128 fine samples for each ray, and render the deformed template image \hat{I}_o and color modification image \hat{I}_m with the same density values. Then, as explained in Sec. 3.3, we use a color compositing layer to obtain the edited view \hat{I} by blending \hat{I}_m into \hat{I}_o , where the color compositing layer is constructed using a compact UNet-like structure (with 2-layer encoder (3 \rightarrow 16 \rightarrow 32) and a symmetrical decoder, all layers comprise 3 \times 3 convolutions). Besides, we can integrate temporal attribute [51, 52] (from 0 to 1) to the input of $F_{\Delta G}$ (with the positional encoding of 4 frequencies), and train the geometric editing on the edited transitions with temporal attributes as conditions, *e.g.*, the dynamic motion effect shown in the supplementary video.

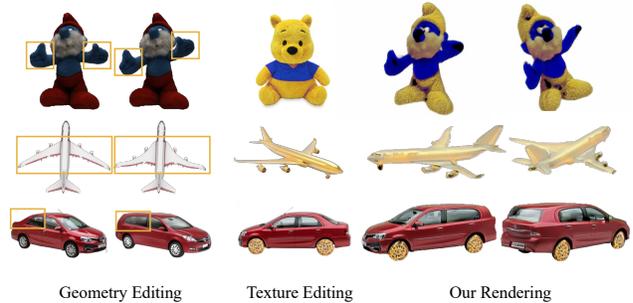


Figure J. We show examples of hybrid object editing by combining geometric and texture editing.

B. Dataset Preparation

We evaluate SINE on both real-world/synthetic and object/scene datasets. Specifically, for the real-world car datasets [57], each sequence contains 72 images with a car rotating on the turntable. We use Colmap [59] to recover camera poses w.r.t the cars' centers for all the images. For the data of Photoshape [50], EditNeRF [37] does not provide edited GT images, so we regenerate all testing cases using Blender, which is more challenging than the original ones (*e.g.*, we stretch the whole chair or enlarge holes, while EditNeRF [37] only fills a tiny hole or removes legs). For the data Blenderswap [39], we render the scenes with Blender's Cycle engine with realistic environment HDR maps. For users' 2D image editing, we use Adobe After Effect / Photoshop to deform images (geometric editing) and paint patterns (texture editing), and use EditGAN [36] and Text2LIVE [1] to edit images with semantic strokes or text-prompts. For users' target images (*e.g.*, cars and chairs) from the Internet, we remove their backgrounds using remove.bg [17] before conducting texture editing.

C. Implementation Details

Training details. As introduced in the main paper, our method performs semantic-driven editing upon the given NeRF model. Specifically, for each object or scene, we first train a generic template NeRF model. Then, we learn geometric editing and texture editing with editing from a single perspective. For geometric editing, since the geometric changes are sometimes combined with minor color changes, we also fine-tune the color modification field with a photometric loss (see the first term in Eq. (3)). For texture editing, the texture transferring loss (Eq.(6)) is defined on a complete image, which is not compatible with NeRF's sparse ray supervision. Therefore, we adopt the deferred

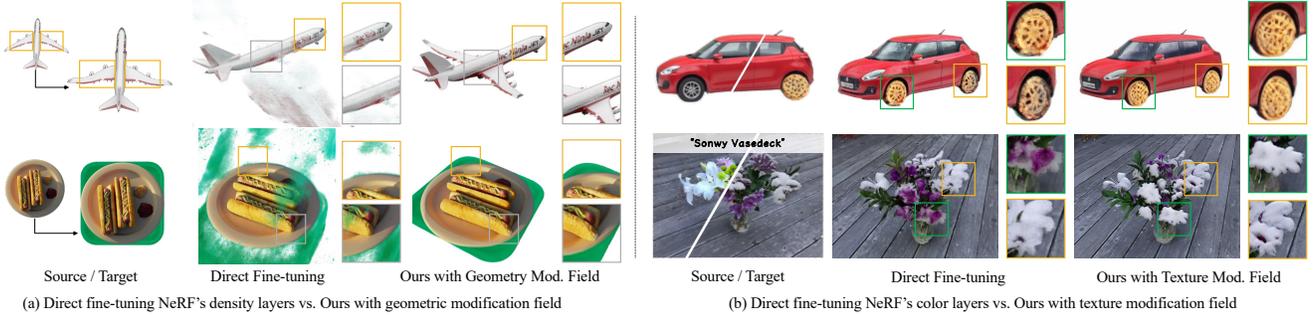


Figure K. We compare our editing field with directly fine-tuning template NeRF.

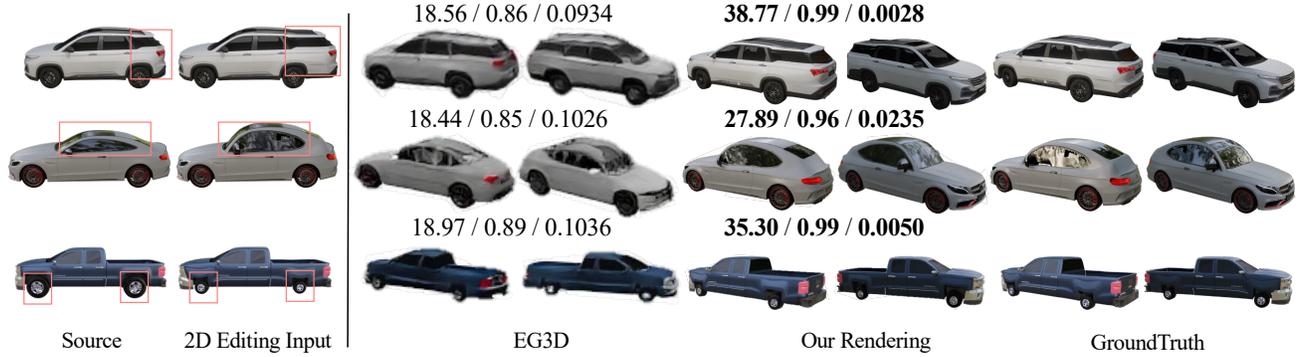


Figure L. We show the quantitative comparison between our method and EG3D [9] on the synthetic cars [39], where the metrics of PSNR \uparrow / SSIM \uparrow / LPIPS \downarrow are annotated above.

back-propagation technique from Zhang *et al.* [79] for texture editing. Practically, we first render the full-sized \hat{I}_o and \hat{I}_m and forward the color compositing layer, and then compute the losses to cache the complete image gradient w.r.t the \hat{I}_m and the color compositing layer, and re-render the \hat{I} and back-propagate the gradients to the $F_{\Delta T}$ and the color compositing layer in a patch-wise manner. To make a smooth convergence, we take the coarse-to-fine regularization [51] on the color modification field by progressively increasing the frequency band of the input features during the training process. Furthermore, we randomly perturb the pose to augment the data distribution and avoid the overfitting of the texture editing. The whole training process of the template NeRF and our editing field takes about 12 hours on a single Nvidia RTX 3090 graphics card.

Preparation of proxy mesh in geometric editing. As introduced in Sec. 3.2, we use a proxy mesh to represent NeRF's geometry during geometric editing. In practice, we directly obtain the proxy mesh using off-the-shelf tools (*i.e.*, implicit surface reconstruction method NeuS [69]). Since we optimize DIF [15] latent code \hat{z} and deform the proxy mesh \hat{M} during the editing, the initial proxy mesh should be binding to a latent code beforehand. Therefore, we obtain the initial latent code \hat{z} to the corresponding initial proxy mesh \hat{M}_σ in an auto-decoding manner [15, 49] before training.

Cycle loss in geometric editing. During the training of

geometric editing, we additionally train a forward modification field $F'_{\Delta G}$ to map the template proxy mesh to the edited space. The forward modification field $F'_{\Delta G}$ and the implicit geometric modification field $F_{\Delta G}$ are both supervised with an cycle loss [35, 43], which is defined as:

$$\mathcal{L}_{\text{cycle}} = \frac{1}{M} \sum_{i=1}^M \|F_{\Delta G}(F'_{\Delta G}(\mathbf{p}_i)) - \mathbf{p}_i\| + \|F'_{\Delta G}(F_{\Delta G}(\mathbf{p}_i)) - \mathbf{p}_i\|, \quad (8)$$

where $\{\mathbf{p}_i | i = 1, \dots, M\}$ are the uniform point samples in 3D space, and we set $M = 1000$ in our experiment.

Feature-cluster-based semantic masking. As introduced in Sec. 3.4, we train a 3D feature field with DINO-ViT's feature maps, and generate feature clusters from the user-painted regions, which will be used to compute semantic masks to distinguish foreground editing areas and background areas. Specifically, we first render the feature map under the specific editing view, and sample 1000 feature points on the user's painted region (which is directly accessible from the editing tools). Then, we use K-Means to generate $K = 15$ clusters from the sampled feature points. During the training stage, we first render the current training view's feature map, and compute the L2-normalized pixel-wise feature distance (from 0 to 1) to the nearest clusters. The pixels with distances smaller than 0.5 would be marked as foreground objects, and the others would be marked as background. These computed editing masks would be used



Figure M. We show the comparison of our method with EditGAN [36] on the real-world car dataset [57].

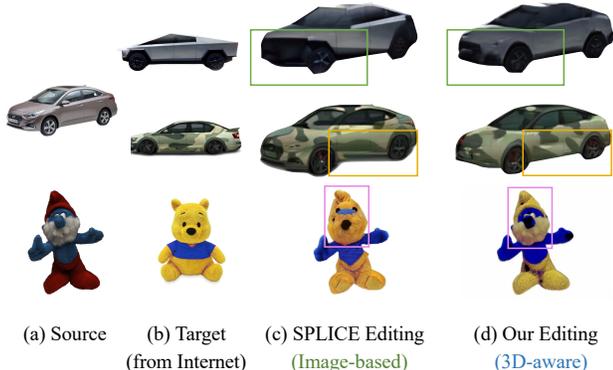


Figure N. We compare our method with SPLICE-ViT [66] on the real-world cars [57] and toys [26]. Since our editing method is built upon 3D-aware models, we consistently achieve better texture-transferring results than SPLICE-ViT when the source and target are observed from different perspectives of views (*e.g.*, cars) or with significant different shapes (*e.g.*, plush toys).

to regularize both geometric and texture editing (see Eq. (7)) to maintain the irrelevant content unchanged.

User study. The questionnaire contains 17 cases, 8 for target-image-based editing (*e.g.*, Fig. 7 (a)) and 9 for text-prompts editing (*e.g.*, Fig. 7 (b)). We show the participants a source image, a target image/text prompts, as well as the results produced by different methods. Participants are asked to select one result that best matches the style of the target image or the text meaning.

D. More Experiments

Hybrid editing with geometric and texture changes. We can combine geometric and texture editing on the same object by optimizing geometric-related losses and texture-transferring losses in turns. As shown in Fig. J, we can edit objects’ geometries while transferring textures with users’ target images, *e.g.*, the plush toy raises its hands and is painted in new textures from a yellow bear, and the airplane extends its wings and is painted golden. Please refer to our supplementary video for a vivid animation of these effects.

3D editing field vs. template NeRF fine-tuning. In this experiment, we compare our 3D editing field with the naïve fine-tuning template NeRF (which is adopted by CLIP-NeRF [68] and DFF [32]). Editing NeRF with only a single image is fairly ambiguous without external supervision (*e.g.*, semantic hints). For a fair comparison, we provide external supervision for the baseline method (vanilla NeRF). Specifically, for texture editing, we enable NeRF’s related network layers to be optimized and use the same texture editing losses. Directly fine-tuning NeRF’s color layers can change the objects’ texture to some extent but cannot reach the same quality as our full model (*e.g.*, uncovered cookie tires and snowy flowers in Fig. K (b)). For geometry editing, we fine-tune vanilla NeRF with \mathcal{L}_{gt} and \mathcal{L}_{reg} since it is not trivial to apply SDF-based shape priors to vanilla NeRF. As demonstrated in Fig. K (a), naively fine-tuning NeRF on geometric editing would lead to the overfitting to a single view, and the multi-view consistency is no longer ensured (*e.g.*, broken wings and green floaters in Fig. K (a)).

Quantitative comparison with EG3D on synthetic cars. We conduct quantitative comparisons with the SOTA 3D-aware GAN method EG3D [9] on the synthetic car dataset. To obtain the ground-truth images of the edited results, we use Blender to render the training and testing views, and modify the cars’ geometry within the software. As shown in Fig L, our method achieves better rendering quality than EG3D on both visual quality and all the metrics (PSNR, SSIM, and LPIPS). For example, we can preserve the specular effect even after the editing (*e.g.*, the specular area facing the light source and the reflection on the windshield), while EG3D struggles to produce photo-realistic results due to the limitation of its learned 3D latent representation.

Comparison with 2D GANs. We compare our method against the SOTA 2D semantic editing method EditGAN [36] on the real-world car dataset [57]. To make a fair comparison, we train our NeRF’s backbone and EditGAN’s style codes on all the multi-view images (*i.e.*, each style code corresponds to one view). Then, we perform semantic 2D editing on one single view using EditGAN. For our method, we use the edited view to train our editing field. And for EditGAN, we save the intermediate editing vector and add the editing vector to all the style codes, which yields multi-view edited images. As demonstrated in Fig M,

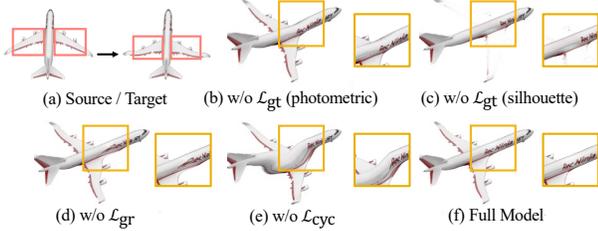


Figure O. We inspect the efficacy of different constraints in geometric editing.



Figure P. We inspect the efficacy of the texture prior constraint in texture editing

since EditGAN is agnostic to the 3D geometry, its results suffer from the inconsistent issue between different views, *e.g.*, poor inversion results for the head and tail of the car, and the semantic editing result cannot be precisely applied to all views. In contrast, our methods can synthesize cars with multi-view consistency and high-quality editing results.

More ablation studies on geometric supervision. As shown in Fig. O, since ablating \mathcal{L}_{gt} (Eq. (3)) makes it no supervision on editing, we split it into photometric (b) and silhouette (c) terms, and the absence of either will result in distorted or washed-out texture. (d): When ablating deformation reg. loss \mathcal{L}_{gr} (Eq. (4)), the edited object is severely distorted (*e.g.*, the letters are stretched). (e): The cycle loss \mathcal{L}_{cyc} (Eq. (1) in supp.) brings constraints from shape prior to geometric mod. field $F_{\Delta G}$, and ablating it would lose the efficacy of semantic guidance (*e.g.*, the twisted airplane).

Ablation study of texture supervision. In Fig. P, we disable texture transfer loss \mathcal{L}_{tex} (Eq. (6)) and utilize photometric loss to paint the target texture, which leads to incomplete texture transferring results for invisible parts as shown below. Besides, without texture prior supervision, we cannot transfer textures between objects with different shapes.

Comparison of texture editing with image-based SPLICE-ViT. Our texture transferring loss (Eq. (6)) is inspired from SPLICE-ViT [66], but fully leverages the multi-view training scheme. Therefore, we compare our texture editing with image-based SPLICE-ViT in Fig. N. As shown in Fig. N, SPLICE-ViT is sensitive to the perspective difference of the source and target images, which results in overfitting appearances on the edited view, *e.g.*, horizontal straight patterns of cars when observing cars from a slightly tilted view, distorted faces of the plush toy. By contrast, our method consistently achieves better texture-transferring results with color patterns properly aligned to the cars' geometries and the plush toy's body parts.

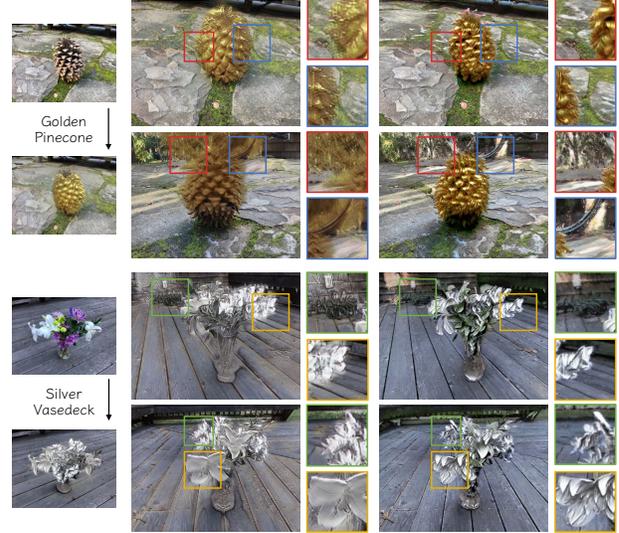


Figure Q. We compare our method with Text2LIVE on texture editing, where our method achieves better multi-view consistency.



Figure R. We show the robustness of geometric editing on proxy meshes with different qualities. The proxy meshes are jittered by adding gaussian noise with different variances (from 0.001^2 to 0.01^2).

Comparison of texture editing with Text2LIVE. As shown in Sec. 4.4 from the main paper, since our method only requires one single-view image as editing input, we can naturally achieve text-prompt-based texture editing by cooperating with off-the-shelf text-driven editing methods (such as Text2LIVE [1]). A follow-up question is, how does the Text2LIVE itself perform to the same 360° dataset in our texture editing task? For video editing, Text2LIVE uses layered atlas [30] to convert objects and backgrounds into separated 2D layers. However, in the unbounded 360° dataset (*e.g.*, pinecone and vase [44]), there is no proper way to unwrap 3D objects and scenes into 2D layers (and we also failed to train layered atlas on these 360° datasets). Therefore, we directly apply its converged editing generator to the multi-view images. As shown in Fig. Q, although Text2LIVE produces similar-looking edited images, it cannot maintain multi-view consistency when the viewpoint changes (*e.g.*, blurry edges at the golden pinecone, uncovered petals at the silver vase, and the occasionally affected background). On the contrary, our method naturally takes advantage of multi-view training and consistently delivers more plausible and realistic novel views.

Robustness to the noise of proxy mesh. The geometry prior guidance uses the proxy mesh to supervise the ge-

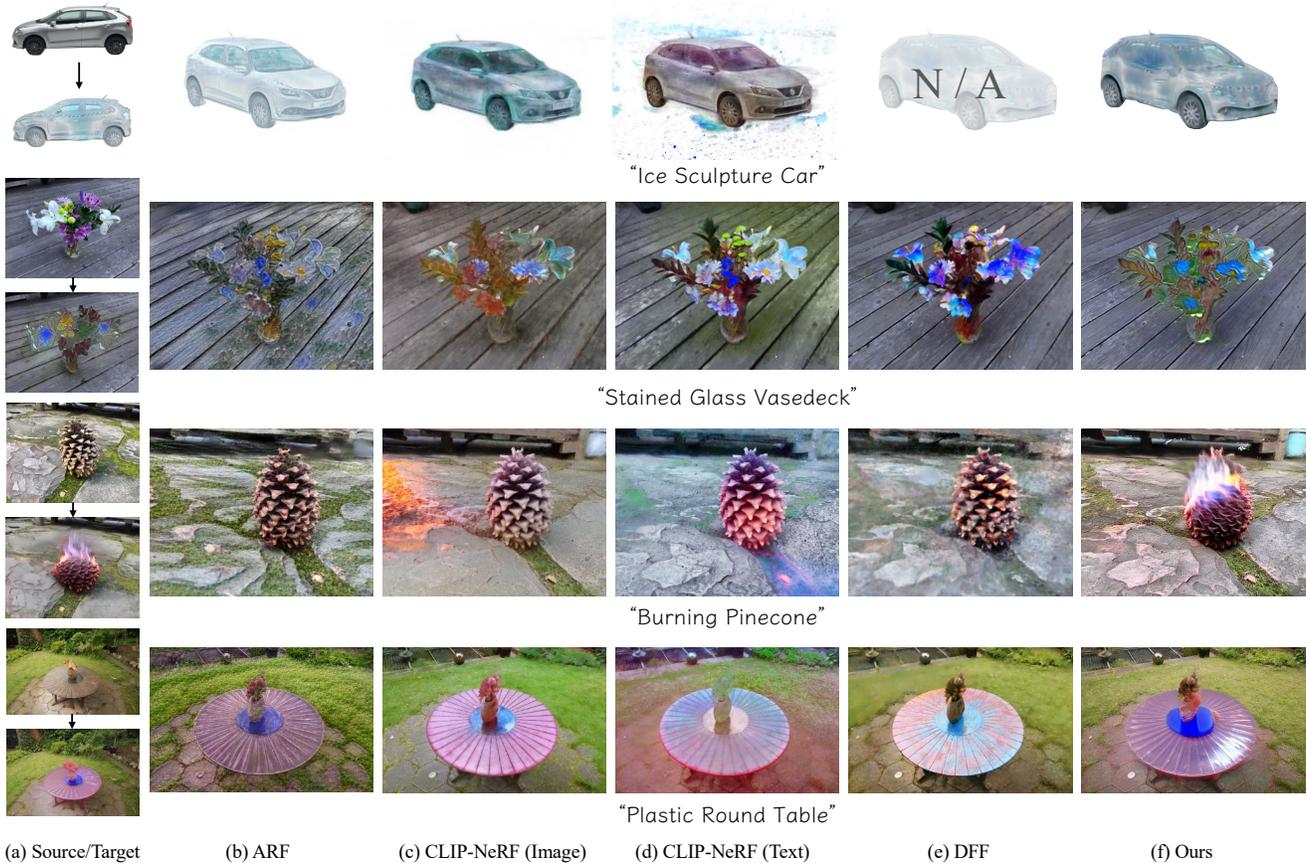


Figure S. We show more comparison results of texture editing with ARF, CLIP-NeRF, and DFF on the real-world car [57] and 360° scene dataset [44]. Our method consistently achieves more realistic and appealing editing results than the others.

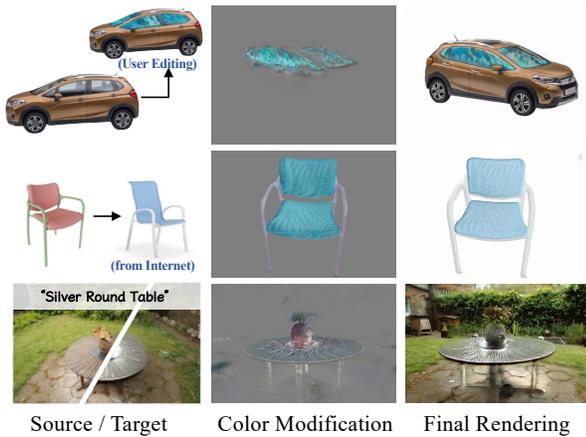


Figure T. We show more rendering results from color modification field and compositional layer.

ometry modification field. Therefore, we analyze the effect of mesh quality on our editing results. Specifically, we add 3 groups of gaussian noise to the vertices of the proxy mesh and conduct training of our editing field. As shown in Fig R, our method can robustly learn geometric editing even with noisy proxy mesh (*e.g.*, with the gaussian noise of $N(0, 0.004^2)$ in the third column).

More comparison results on texture editing. We show more comparison results on the texture editing task with ARF [79], CLIP-NeRF [68] and DFF [32] in Fig. S (where Fig. S (a) is the source view and the target view produced by Text2LIVE [1]). For CLIP-NeRF [68], since the official codebase has not been fully released, we use our own implementation by fine-tuning NeRF’s color-related field with CLIP loss, and both use the target features from text embedding and the image embedding (with the same target images in Fig. S (a)), which are denoted as CLIP-NeRF (Image) and CLIP-NeRF (Text), respectively. For DFF [32], we adopt the official codebase and use the texts for NeRF editing and background ray-filtering according to the document. In Fig. S, we omit the DFF’s object-centric comparison on the car, since it mainly focuses on scene-level decomposition and editing. As demonstrated in Fig. S, NeRF stylization methods like ARF cannot precisely edit fine-grained effects on the desired location. NeRF fine-tuning approaches like CLIP-NeRF and DFF only change appearance colors, but cannot produce vivid effects (*e.g.*, the burning pinecone or ice sculpture cars). Note that although DFF uses the semantic-field guided decomposed rendering to maintain the background color unchanged, this strategy is not com-



Figure U. We show the results of geometry editing with topology changes.

patible with our color compositing mechanism since we introduce an additional 2D CNN layer to blend the template and editing color for better visual appearance.

By contrast, our method both achieves realistic and appealing editing effects, and also effectively preserves background content, and the results are consistently preferred by most of the participants in the user study (see Sec. 4.4).

Impact of texture modification field & color compositional layer. The texture modification field learns detailed modifications and the compositional layer blends the original and modified rendering to produce the final edited results, as demonstrated in Sec. 4.5 and Fig. 8 (a). Here we show more rendered texture mod. field (a.k.a. color mod. \hat{I}_m) in Fig T.

Deformation with topology changes. Our method does not support deformation with topology changes such as breaking the plate, but can provide a visually plausible result by making the “broken part” white, as shown in Fig U. In the future, we can integrate more flexible representations such as ambient slicing surface [52] into our model.