

NeuralEditor: Editing Neural Radiance Fields via Manipulating Point Clouds

Jun-Kun Chen^{1†} Jipeng Lyu^{2†} Yu-Xiong Wang¹

¹University of Illinois at Urbana-Champaign ²Peking University [†]Equal Contribution

{junkun3, yxw}@illinois.edu lvjipeng@pku.edu.cn

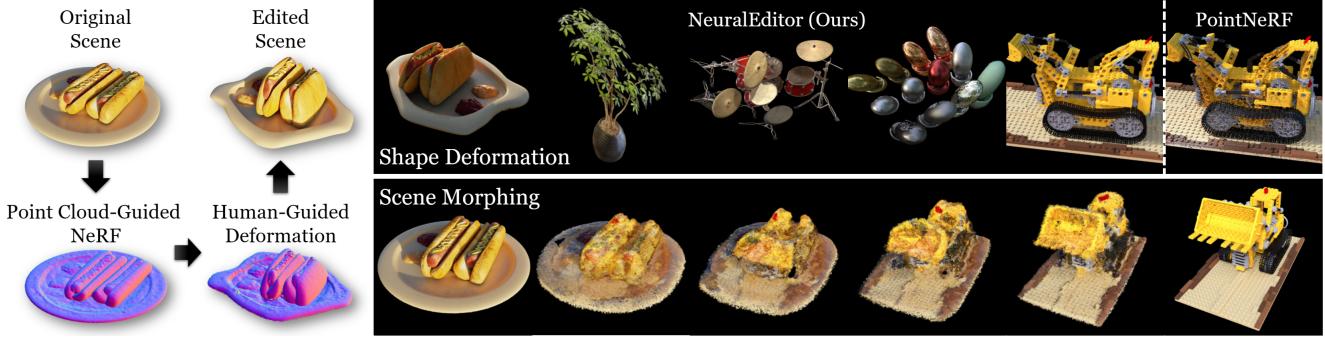


Figure 1. Our NeuralEditor offers *native* support for general and flexible shape editing of neural radiance fields via manipulating point clouds. By generating a precise point cloud of the scene with a novel point cloud-guided NeRF model, our NeuralEditor produces high-fidelity rendering results in both shape deformation and more challenging scene morphing tasks.

Abstract

This paper proposes *NeuralEditor* that enables neural radiance fields (NeRFs) natively editable for general shape editing tasks. Despite their impressive results on novel-view synthesis, it remains a fundamental challenge for NeRFs to edit the shape of the scene. Our key insight is to exploit the explicit point cloud representation as the underlying structure to construct NeRFs, inspired by the intuitive interpretation of NeRF rendering as a process that projects or “plots” the associated 3D point cloud to a 2D image plane. To this end, *NeuralEditor* introduces a novel rendering scheme based on deterministic integration within K-D tree-guided density-adaptive voxels, which produces both high-quality rendering results and precise point clouds through optimization. *NeuralEditor* then performs shape editing via mapping associated points between point clouds. Extensive evaluation shows that *NeuralEditor* achieves state-of-the-art performance in both shape deformation and scene morphing tasks. Notably, *NeuralEditor* supports both zero-shot inference and further fine-tuning over the edited scene. Our code, benchmark, and demo video are available at immortalco.github.io/NeuralEditor.

1. Introduction

Perhaps the most memorable shot of the film *Transformers*, *Optimus Prime* is seamlessly transformed between a

humanoid and a Peterbilt truck – such free-form editing of 3D objects and scenes is a fundamental task in 3D computer vision and computer graphics, directly impacting applications such as visual simulation, movie, and game industries. In these applications, often we are required to manipulate a scene or objects in the scene by editing or modifying its shape, color, lighting condition, *etc.*, and generate visually-faithful rendering results on the edited scene efficiently. Among the various editing operations, shape editing has received continued attention but remains challenging, where the scene is deformed in a human-guided way, while all of its visual attributes (*e.g.*, shape, color, brightness, and lighting condition) are supposed to be natural and consistent with the ambient environment.

State-of-the-art rendering models are based on implicit neural representations, as exemplified by neural radiance field (NeRF) [27] and its variants [3, 33, 37, 39, 48]. Despite their impressive novel-view synthesis results, most of the NeRF models substantially lack the ability for users to adjust, edit, or modify the shape of scene objects. On the other hand, shape editing operations can be natively applied to explicit 3D representations such as point clouds and meshes.

Inspired by this, we propose *NeuralEditor* – a general and flexible approach to editing neural radiance fields via manipulating point clouds (Fig. 1). Our *key insight* is to benefit from the best of both worlds: the superiority in rendering performance from implicit neural representation

combined with the ease of editing from explicit point cloud representation. NeuralEditor enables us to perform a wide spectrum of shape editing operations in a consistent way.

Such introduction of point clouds into NeRF for general shape editing is rooted in our interpretation of *NeRF rendering as a process that projects or “plots” the associated 3D point cloud to a 2D image plane*. Conceptually, with a dense enough point cloud where each point has an opacity and its color is defined as a function of viewing direction, directly plotting the point cloud would achieve similar visual effects (*i.e.*, transparency and view-dependent colors) that are rendered by NeRF. This *intrinsic integration* between NeRF and point clouds underscores the advantage of our NeuralEditor over existing mesh-based NeRF editing methods such as NeRF-Editing [51], Deforming-NeRF [44], and CageNeRF [30], where the process of constructing and optimizing the mesh is separated from the NeRF modeling, making them time-consuming. More importantly, with the point cloud constructed for a scene, the shape editing can be natively defined as and easily solved by just moving each point into the new, edited position and re-plotting the point cloud. Therefore, our approach supports more general scene editing operations which are difficult to achieve via mesh-guided space deformation.

The key component in our NeuralEditor lies in a point cloud-guided NeRF model that natively supports general shape editing operations. While the recent method PointNeRF [43] has demonstrated improved novel-view synthesis capability based on point clouds, it is not supportive to shape editing. Our idea then is to exploit the underlying point cloud in ways of not only optimizing its structure and features (*e.g.*, adaptive voxels) for rendering, but also extracting additional useful attributes (*e.g.*, normal vectors) to guide the editing process. To this end, we introduce K-D trees [4] to construct density-adaptive voxels for efficient and stable rendering, together with a novel deterministic integration strategy. Moreover, we model the color with the Phong reflection [31] to decompose the specular color and better represent the scene geometry.

With a much more precise point cloud attributed to these improvements, our NeuralEditor achieves high-fidelity rendering results on deformed scenes compared with prior work as shown in Fig. 1, even in a *zero-shot* inference manner without additional training. Through fast fine-tuning, the visual quality of the deformed scene is further enhanced, almost perfectly consistent with the surrounding lighting condition. In addition, under the guidance of a point cloud diffusion model [24], NeuralEditor can be naturally extended for smooth *scene morphing* across multiple scenes, which is difficult for existing NeRF editing work.

Our contributions are four-fold. (1) We introduce NeuralEditor, a flexible and versatile approach that makes neural radiance fields editable through manipulating point

clouds. (2) We propose a point cloud-guided NeRF model based on K-D trees and deterministic integration, which produces precise point clouds and supports general scene editing. (3) Due to the lack of publicly available benchmarks for shape editing, we construct and release a reproducible benchmark that promotes future research on shape editing. (4) We investigate a wide range of shape editing tasks, covering both shape deformation (as studied in existing NeRF editing work) and challenging scene morphing (a novel task addressed here). NeuralEditor achieves state-of-the-art performance on all shape editing tasks in a unified framework, without extra information or supervision.

2. Related Work

Neural Scene Representation. Traditional methods model scenes with explicit [2, 11, 15, 20, 34, 38] or implicit [6, 19, 26, 28, 40] 3D geometric or shape representations. Initiated by NeRF [27], leveraging implicit neural networks to represent scenes and perform novel-view synthesis has become a fast-developing field in 3D vision [8, 9]. While most of the follow-up work focuses on improving aspects such as the rendering realism [3, 10, 21, 37], efficiency [33, 41, 48], and cross-scene generalization [5, 39, 43, 49], the scene editing capability is substantially missing in the NeRF family which we address in this paper. In addition, we exploit K-D tree-guided point clouds as the underlying structure, different from other NeRF variants based on octrees [21, 33, 48] or plain voxels [41].

Point-Based NeRFs. Recently, using point clouds to build a NeRF model has shown better encoding of scene shape and improved rendering performance, as represented by PointNeRF [43]. PointNeRF proposes a point initialization network to produce the initial point cloud together with the point features, which is further optimized by a pruning and growing strategy. While both PointNeRF and our NeuralEditor employ point clouds as the underlying structure, NeuralEditor better exploits useful information within the point clouds: PointNeRF only directly uses the locations of points; by contrast, NeuralEditor *considers the point cloud more as a geometrical shape* and extracts relevant information like normal vectors, which plays an important role in rendering and shape editing. Importantly, our approach is designed to support scene editing, in contrast to PointNeRF.

Scene Editing via NeRFs. Different types of scene editing have been studied under NeRFs. EditNeRF [22], ObjectNeRF [45], and DistillNeRF [18] perform simple shape and color editing for objects specified with human-input scribble, pixel, segment, language, *etc.* Neu-Physics [32] edits a dynamic scene via physics parameters. CCNeRF [35] proposes an explicit NeRF representation with tensor rank decomposition to support scene composition. INSP-Net [42] considers filter editing like denoising. Such work cannot address 3D shape editing and only sup-

ports simple editing operations, like object selection, similarity transformation, or limited shape deformation.

3D Shape Editing. Traditional representation methods support keypoint-based shape editing [13, 14, 16, 25, 36, 47, 52] with meshes [29, 50], which cannot be directly applied to implicit representations used by NeRF. Existing NeRF editing work primarily studies a particular shape editing task, mesh deformation, and addresses it in a common paradigm [30, 44, 51]: A mesh of the scene is first constructed by either exporting it from a trained NeRF with the Marching Cubes algorithm [23], or optimizing close-to-surface cages along with training. After the user deforms the mesh, the deformed scene is rendered by deforming the space and bending the viewing rays in the original scene with the trained NeRF. Doing so requires extra efforts to convert *implicit* scene representation to *explicit* mesh, which might not be precise enough, and only supports *continuous* shape editing that can be converted to space deformation. On the contrary, our NeuralEditor directly maintains and utilizes alternative explicit scene representation – *the point cloud which is intrinsically integrated with NeRF*, making NeuralEditor require no extra efforts and support more general shape editing tasks like scene morphing. NeuralEditor supports both zero-shot inference and further fine-tuning over the edited scene, while prior work cannot.

3. NeuralEditor: Point Cloud-Guided NeRF

We propose a novel point cloud-guided NeRF model, NeuralEditor – it not only achieves realistic rendering results in the novel-view synthesis task, but also produces a point cloud that precisely describes the shape of the scene, thus facilitating general shape editing tasks. As illustrated in Fig. 2, we leverage the K-D trees [4] to construct *density-adaptive* voxels (which also naturally enable us to skip empty spaces), and introduce *deterministic* spline integration for rendering. We use the Phong reflection to model the color along with the normal vectors obtained from the underlying point cloud. With our enhanced point cloud optimization, NeuralEditor obtains much more precise underlying point clouds, compared with noisy and imprecise outputs of state-of-the-art PointNeRF [43] (as shown in Sec. 5).

3.1. K-D Tree-Guided Voxels

To render with points, we construct multi-scale density-adaptive voxels based on K-D trees [4], namely, *K-D voxels*. K-D trees are a data structure constructed on K -dimensional points, where $K = 3$ for 3D points. As a special decision tree, K-D tree’s each node divides the point set into two equal-sized parts with axis-parallel criterion.

For each K-D tree’s node, we compute its bounding box by taking the minimum and maximum x, y, z coordinates in its subtree and with proper padding margins. As we divide the points in a top-down manner in one of the x, y, z direc-

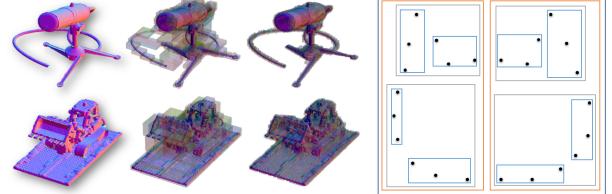


Figure 2. **Our K-D Voxels.** **Column 1:** original point clouds constructed from two scenes, colored with normal vector directions. **Column 2:** upper-level voxels which coarsely represent the shape. **Column 3:** lower-level voxels which tightly cover the shape. **Right:** Visualization of K-D voxels on a 2D point cloud. Each color represents boxes of nodes on each level of the K-D tree. Lower-level boxes containing fewer points cover the shape more tightly, and vice versa for higher-level boxes.

tions, in each layer of the K-D tree, different nodes’ bounding boxes are *mutually exclusive*. Therefore, the bounding boxes can be regarded as voxels. As boxes in the upper layers contain more points (larger voxels), while those in the lower layers contain fewer points (smaller voxels), we natively obtain a *multi-scale* voxel construction from one K-D tree. As shown in Fig. 2, voxels from the large to small scales represent the shape of the scene from coarse to fine.

3.2. Rendering Over K-D Voxels

We now introduce a rendering scheme that exploits K-D voxels to perform all the sub-procedures associated with rendering in a unified way. This scheme enables us to render more naturally, efficiently, and even deterministically, meanwhile it also simplifies some widely-adopted design choices in conventional NeRF rendering.

Skipping Empty Spaces. In NeRF rendering, we are supposed to focus only on the *surface* of scene objects. As shown in Fig. 2, all our K-D voxels are produced to stick to the surface of objects, which the point cloud is constructed to describe. Such a property allows us to avoid explicitly “skipping” empty spaces, which often requires extra consideration in most NeRF models – only considering the space inside a voxel automatically focuses on the surface; as the depth of the voxel’s node goes deeper, it becomes closer to the surface. Moreover, during the construction of the K-D tree, the points at each node are divided within its sub-nodes, and the node’s voxel fully covers all its sub-nodes. This further provides us with a native *top-down recursive* procedure to locate the voxels intersected with the querying ray: We start from the root node, and recurse on the sub-nodes until we (1) reach a pre-set node depth (or equivalently, a pre-set voxel size) and then query within the associated voxel, or (2) stop recursion on non-intersected nodes.

Density-Adaptive Rendering. An important design in NeRFs is the coarse-to-fine strategy for density-adaptive rendering, so that more points are sampled in high-volume density areas. Our K-D voxels natively support such a design *without additional bells and whistles*. This is because

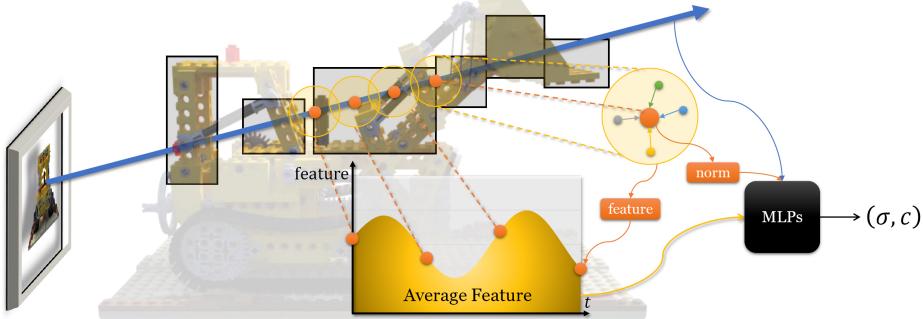


Figure 3. **Our NeuralEditor architecture.** We propose deterministic spline integration for KNN-based point features over each K-D tree-guided density-adaptive voxel, and model the color via Phong reflection with normal vectors estimated from the point cloud’s shape.

voxels in the same K-D tree layer contain the same number of points. As the point density can be regarded as an approximation of volume density, all such voxels have the same density. Therefore, we directly use K-D voxels to guide the density-adaptive rendering. Specifically, we conduct the rendering process at the voxels of some bottom layers in the K-D tree. For each querying ray, we use the aforementioned recursive procedure to locate the minimal intersected voxels that are deep enough. Here “minimal” means that the ray intersects with the node’s voxel, but does not intersect with any sub-node’s voxel. These intersected voxels divide the querying ray into several segments (Fig. 3). The segments covered by a voxel are close to the surface and used for rendering, while those not covered are in empty spaces.

Deterministic Spline Integration. DIVeR [41] shows that deterministic integration outperforms stochastic integration in NeRF rendering. So we perform a deterministic integration to obtain the segment’s feature within each voxel. Since we do not necessarily have points at the voxel’s vertices, the trilinear interpolation used in DIVeR is not feasible here. Instead, we use spline integration. For the i -th intersected voxel in the ray passing order, we uniformly select points in this segment, and integrate their features to obtain the average feature f_i of the segment of the i -th voxel:

$$f_i = \frac{1}{r_i - l_i} \int_{l_i}^{r_i} \mathbf{feature}(o + t \cdot d) dt, \quad (1)$$

where $[l_i, r_i]$ is the intersection interval, and o and d are the source and direction of the querying ray, respectively. This average feature f_i can be interpreted as the feature of a *representative point* p_i located somewhere in the segment.

KNN-Based Feature Aggregation. For each uniformly selected point q in the segment during spline integration, we obtain its feature via weighted interpolation from the features of its K nearest neighbors (KNN) in the point cloud:

$$\begin{aligned} \mathbf{feature}(q) &= \sum_{p_j \in \text{KNN}(q; K)} k_j e_j, \\ \{k_j\} &= \underset{p_j \in \text{KNN}(q; K)}{\text{SoftMax}} \left(\log \gamma_j - \log \|q - p_j\|_2^2 \right), \end{aligned} \quad (2)$$

where for each point p_j in the point cloud, we parameterize its confidence γ_j and point feature e_j as in PointNeRF.

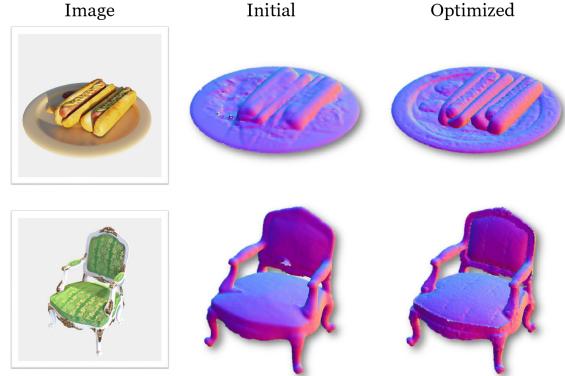


Figure 4. In the two scenes of NeRF Synthetic [27], NeuralEditor optimizes the rough initial point cloud to a precise point cloud. The points are colored with their normal vectors.

Phong Reflection-Based Color Modeling with Point Cloud Normal Vectors. To obtain the volume density σ_i and color c_i of the representative point p_i in the i -th voxel, we use the Phong reflection model [31]. As we have the underlying point cloud, we use Open3D [53] to estimate the normal vector for each point, and integrate these vectors over the interval to get an average normal vector n_i . Such information better characterizes the shape of point clouds (scenes), which plays an important role in Phong-based rendering and also implicitly facilitates optimizing the point clouds (Sec. 3.3). Consistent with RefNeRF [37], we use multiple multilayer perceptrons (MLPs) to model other attributes, including volume density, tint, roughness, and diffuse and specular color, and use the Phong formula to calculate the final c_i with these attributes. Finally, we aggregate c_i of all segments on the ray to obtain the final color:

$$\begin{aligned} c_{\text{pixel}} &= \sum_{i \geq 1} \tau_i \cdot (1 - \exp(-(r_i - l_i)\sigma_i)) \cdot c_i, \\ \tau_i &= \exp \left(- \sum_{i'=1}^{i-1} (r_{i'} - l_{i'})\sigma_{i'} \right). \end{aligned} \quad (3)$$

3.3. Point Cloud Optimization

Point Cloud Initialization. To start training, we need a coarse initial point cloud. Consistent with PointNeRF, we

use a point generation network, which consists of a multi-view stereo (MVS) model [12] based on a 3D convolutional neural network (CNN), to generate the points' coordinates and confidence values, and another 2D CNN [46] to generate their initial features. This network was pre-trained on the DTU training dataset [1], and can generalize to other datasets and scenes. As shown in Fig. 4, the initial point cloud generated by such a network is coarse and noisy.

Explicit Optimization via Pruning and Growing. We perform a similar pruning and growing procedure as in PointNeRF, to prune outliers with low confidence γ_j and fill holes in the point clouds. We make several important modifications over PointNeRF, and integrate this procedure with our deterministic integration (details in the supplementary).

Implicit Optimization with Normal Vectors. In addition to the explicit optimization, the point cloud is also optimized implicitly during training through the adjustment of point confidence γ_j . When computing the average normal vector for rendering, we aggregate normal vectors of nearby points weighted with their distance and confidence, where the confidence of noisy or inaccurate points with potentially abnormal normal vectors is adjusted accordingly. Moreover, we apply the normal vector regularization losses from RefNeRF [37] to supervise the points' confidence w.r.t. their normal vectors. These strategies collectively provide *implicit but more tailored* ways to optimize the point clouds. With both explicit and implicit optimization, NeuralEditor obtains very precise point clouds (Fig. 4).

4. Shape Editing with NeuralEditor

Formulation of General Shape Editing Tasks. We define the shape editing tasks based on *indexed* point clouds. To this end, we first re-define an indexed point cloud P as a mapping from a point index j to the corresponding point p_j ,

$$P : j \rightarrow p_j, \text{ where } p_j \in \mathbb{R}^3, j = 1, \dots, |P|. \quad (4)$$

A shape editing task is defined as *another* indexed point cloud $Q(P)$,

$$Q(P) : j \rightarrow q_j, \text{ where } q_j \in \mathbb{R}^3 \cup \{\emptyset\}, j = 1, \dots, |P|, \quad (5)$$

describing a shape editing task whether the j -th point moves from p_j to q_j or is deleted in the deformation if $q_j = \emptyset$. With this definition, $Q(P)$ can be an *arbitrary* point cloud with points properly matched to points in P by same indices, regardless of connectivity or continuity.

Our formulation represents a broad range of shape editing tasks. The mesh editing tasks in NeRF-Editing [51], Deforming-NeRF [44], and CageNeRF [30] can be more simply and clearly defined here. For example, in NeRF-Editing, a mesh is exported from a general NeRF model, deformed manually, and converted to an “offset” or a continuous space deformation. We can depict such a task without

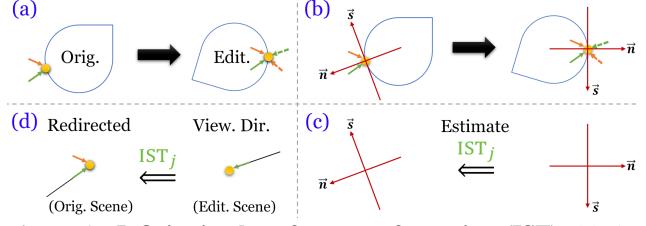


Figure 5. **Infinitesimal surface transformation (IST).** (a) As the view-dependent colors are modeled as absolute viewing directions, they (solid arrows at the right) are different from the correct colors (dashed arrows at the right) after deformation. We solve this by (b) constructing a local coordinate system near the j -th point and (c) modeling IST from the edited scene to the original scene with the coordinate systems, so as to (d) redirect the viewing direction to the original scene when rendering the edited scene.

“offsets,” by recording only the *final location for each point* without extra information. Notably, our formulation even models those whose deformation is *not continuous* in the space, e.g., cutting a scene into two parts, and thus cannot be covered and solved by NeRF-Editing, Deforming-NeRF, or CageNeRF, as shown in the supplementary.

Editing Shape by Moving Points. We design our shape editing scheme with NeuralEditor. This is achieved by interpreting NeRF rendering as “plotting” the sampled points over the viewing ray. If we render the scene by naively plotting the point cloud P , the shape editing task can be addressed by replacing each point’s coordinate from p_j to q_j . For NeuralEditor, we similarly replace the underlying point cloud from P to $Q(P)$, while maintaining the confidence values and features. This method is general and can also be applied to any point-based NeRF model like PointNeRF.

Correcting View-Dependence with Infinitesimal Surface Transformation (IST). The editing method above can already obtain reasonable results. However, as illustrated in Fig. 5, the modeled view-dependent colors record the *absolute* viewing direction, making them incorrect after deformations that change their orientation.

To solve this issue, we model the infinitesimal surface transformation (IST) for each point to redirect the viewing ray in the correct direction. We construct a local coordinate system for each point to represent the orientation of the infinitesimal surface, using its normal vector and two point indices that are neighbors of the j -th point in both P and $Q(P)$. By comparing these two coordinate systems, we can obtain an affine transformation IST_j for the j -th point to redirect the querying view direction (Fig. 5). This procedure is different from modeling space deformation [30, 44, 51], as we only need to model a simple affine transformation at each point, while those methods model a complicated, continuous, and non-linear deformation in the whole space.

Our proposed method requires a precise point cloud with normal vector-based color modeling. It is thus incompatible with PointNeRF, as PointNeRF is unable to obtain a desired point cloud to estimate the surface normal vectors.

Fine-Tuning on Deformed Scene. Using the shape editing scheme introduced above, we can apply shape deformation on the scene modeled by our NeuralEditor *without any modification* to the model architecture or rendering scheme, which means that the resulting model is still a valid, fully functional NeuralEditor. Therefore, we can further fine-tune NeuralEditor on the deformed scene if the ground truth is available. We can even fine-tune the infinitesimal surface transformation with other parameters, to rapidly adjust toward better ambient consistency. This makes NeuralEditor desirable *in practice*, since in most applications, the final goal is not a zero-shot inference, but to *fit* the deformed scene well with reduced cost. By supporting fine-tuning, our NeuralEditor aligns well with and achieves this goal.

As another point-based NeRF model, PointNeRF supports fine-tuning but cannot leverage infinitesimal surface transformation fine-tuning to further optimize the performance. On the other hand, mesh-based NeRF editing models [30, 44, 51] do not support fine-tuning well: With deforming the space instead of the scene, these models’ rendering scheme has highly changed. In their rendering process, a ray may go through a long, irregular way to reach the scene’s surface. As the modeled space deformation might not be precise, it could be hard to tune the irregular space well, and even hurt other parts of the trained NeRF model. Such issues occur especially for some spaces with non-uniform density, since most of their model components (*e.g.*, positional encodings, voxels) are not designed to deal with non-uniform spaces. Among all these methods, only our NeuralEditor has complete support for fine-tuning.

5. Experiment

Point Cloud Generation. The underlying point cloud is fundamental to all editing tasks. Fig. 6 first provides a qualitative comparison of point clouds generated by our NeuralEditor and PointNeRF [43] on NeRF Synthetic [27]. Ours are much more precise with sharper details, *e.g.*, the mayonnaise on the Hotdog’s sausage, the uneven texture on the Chair’s cushion, the edge of the Mic’s stand, and the Lego brick’s studs. By contrast, PointNeRF’s point clouds are blurred and noisy, lose most of the details, and even contain obvious shape defects on the Hotdog’s plate and Chair’s backrest. This shows that while the point cloud generation task is challenging, NeuralEditor generates a super-precise point cloud which is crucial for shape editing tasks.

Experimental Settings. We mainly conduct experiments based on scenes from the NeRF Synthetic (NS) dataset. NS is a widely-used NeRF benchmark constructed from Blender [7] scenes. Due to the lack of publicly available benchmarks for shape editing, we use Blender to construct a reproducible benchmark, including the ground truth of edited scenes for evaluation and fine-tuning. Our shape editing tasks cover *all eight scenes* in NS, while prior

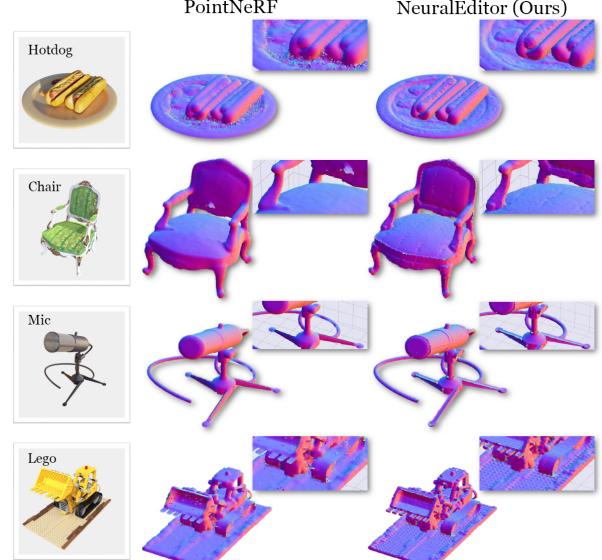


Figure 6. NeuralEditor generates much more precise point clouds than PointNeRF [43] in the four scenes of NeRF Synthetic [27]. The points are colored with their normal vectors.



Figure 7. With too coarse cages, DeformingNeRF [44] is unable to perform the deformation faithfully, leading to poor results.

work [30, 44, 51] only picks a few scenes. The provided images for NS scenes are with opacity, and there is no requirement for the background color. We evaluate and visualize the results on a black background, for better contrast and clearer detail visualization. In the supplementary, we show the results on a white background with same conclusions.

Shape Editing Tasks. We evaluate our model on two types of shape editing tasks, as shown in Fig. 1:

(I) *Shape (Mesh) Deformation Task.* We consider the shape deformation task as in [30, 44, 51]: deform the shape of a scene in a human-guided way. To construct our deformation tasks from NS and obtain the ground truth, we apply the shape deformation simultaneously to the scene and our point cloud within the provided Blender file. We perform both zero-shot inference and fine-tuning, and compare our rendering results with the ground truth. As demonstrated in Figs. 7 and 8, our deformation tasks are much more precise and *aggressive*, compared with those in previous work [30, 44, 51]. In the supplementary, we also design a non-continuous deformation task and deformation tasks on the real-world dataset Tanks and Temples [17] (with zero-shot inference only, as there is no ground truth available).

(II) *Scene Morphing Task.* We address a more challenging shape editing task that has *not* been investigated in prior NeRF editing work: the scene morphing task. Given two scenes *A* and *B*, we should construct a path to gradually

Model	Zero-Shot Inference, PSNR ↑								Fine-Tune for 1 Epoch, PSNR ↑							
	Chair	Hotdog	Lego	Drums	Ficus	Materials	Mic	Ship	Chair	Hotdog	Lego	Drums	Ficus	Materials	Mic	Ship
DeformingNeRF [44]	18.84	-	13.10	-	-	-	-	-	30.11	36.08	31.45	27.16	31.48	27.55	34.34	28.90
PointNeRF [43]	22.21	25.95	24.56	21.00	24.24	21.21	26.77	21.19	32.01	36.38	31.72	28.09	33.21	30.31	35.15	30.01
Naive Plotting	24.91	27.01	25.64	21.29	26.22	21.65	27.63	22.29	32.24	36.69	32.79	28.30	33.34	30.40	35.28	30.08
NeuralEditor w/o IST	24.92	27.02	25.65	21.29	26.24	21.64	27.64	22.28	32.53	37.22	32.95	28.35	33.53	30.82	35.46	30.44
NeuralEditor (Ours)	25.85	27.49	27.46	21.84	27.19	23.18	27.75	24.16	32.53	37.22	32.95	28.35	33.53	30.82	35.46	30.44

Table 1. NeuralEditor significantly and consistently outperforms PointNeRF and Naive Plotting on all deformed scenes of NeRF Synthetic in peak signal-to-noise ratio (PSNR), in both zero-shot inference and fine-tuning settings. Our infinitesimal surface transformation (IST) effectively improves the results by correcting the view-dependent colors. With the precise point cloud generated by NeuralEditor, even Naive Plotting consistently outperforms PointNeRF. Comparison results under other metrics are in the supplementary.

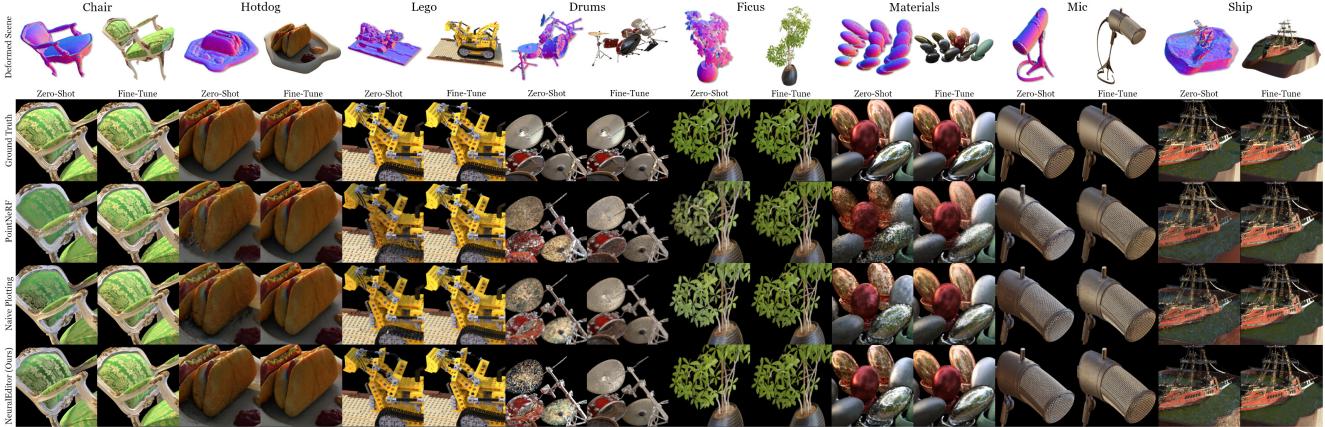


Figure 8. NeuralEditor produces superior rendering results to PointNeRF, with significantly fewer artifacts in zero-shot inference. Fine-tuning further improves the *consistency of rendering with the ambient environment*. We use a black background for better visualization.

change one scene to the other, and the intermediate scenes should have reasonable appearances. We are required to render all intermediate scenes. For this task, we use the point cloud diffusion model [24] to generate intermediate point clouds with *latent space interpolation* between A and B as in [24], and we introduce a K-D tree-based [4] matching algorithm to match the adjacent points to fix the indices of the intermediate scenes. To render an intermediate scene, we apply the shape transformation to the NeuralEditor models trained for scenes A and B , and then interpolate the rendering features to obtain the features of the intermediate scene for rendering.

NeuralEditor Variants. (1) Full NeuralEditor: Our complete model with all components. (2) NeuralEditor *without* infinitesimal surface transformation (IST): We remove the maintenance and optimization of infinitesimal surface transformation in scene editing. This key variant of NeuralEditor enables us to evaluate the importance of IST as well as other components of our model. The *full ablation study* of NeuralEditor is in the supplementary.

Baselines. We compare NeuralEditor against different types of baselines as follows. (1) Naive Plotting: We use the point cloud generated by NeuralEditor, computing each point’s opacity and view-dependent colors with their point features. We render the scene by directly plotting/projecting the point cloud to the camera plane. (2) PointNeRF [43]: For the shape deformation task, we apply the same deformation to the point clouds generated by PointNeRF. For the scene morphing task, we apply the matching algorithm to

the point clouds generated by PointNeRF and the same intermediate point clouds generated by the point cloud diffusion model [24] for fairness. (3) DeformingNeRF [44]: DeformingNeRF is not compatible with the scene morphing task. For the shape deformation task, we perform the same deformation on vertices of given cages. Note that DeformingNeRF only released trained models for Lego and Chair, so we can only evaluate it on these two scenes. While other models [30, 51] support NeRF-based shape deformation via cages or exported meshes, we were unable to use them as baselines – they did not provide executable code nor their deformed scenes for us to evaluate on their tasks.

Shape (Mesh) Deformation Results. The qualitative comparison is shown in Fig. 8. Both PointNeRF and Naive Plotting have many artifacts, like blur, wrong color, black or white shadows, noise, etc., whereas our powerful NeuralEditor produces clean and realistic rendering results. After fine-tuning, NeuralEditor shows a significant improvement with better rendering results than PointNeRF, indicating that NeuralEditor is able to achieve higher consistency with the surrounding ambient environment. Notably, in the Materials scene (the 6th scene from left), only our NeuralEditor generates reasonable reflection, while both baselines show blurry and visually messy results. Also in the Drums scene, the bottom face of the gong is not visible in any of the training views in the original scene, so all models render poor results in zero-shot inference. However, after only fast fine-tuning, NeuralEditor is able to precisely model the previously unknown surface and generate a sub-

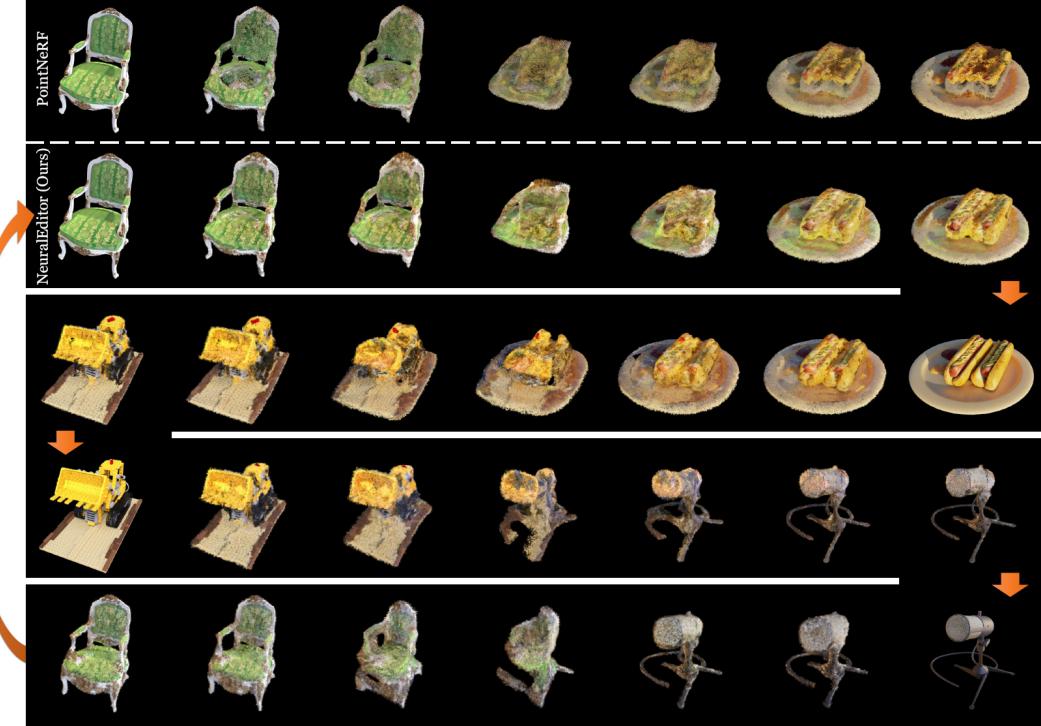


Figure 9. Our NeuralEditor produces *smooth morphing* results between Chair, Hotdog, Lego, and Mic in the NeRF Synthetic dataset, while PointNeRF produces results with blurry textures, black shadows, and gloomy, non-smooth colors. The rendering results in the looped morphing process are arranged in the shape of the numerical digit “3,” indicated by the dividing lines and arrows.

stantially better result than PointNeRF, highlighting NeuralEditor’s strength in fast-fitting. All these results demonstrate that NeuralEditor can handle various visual effects and make them consistent in the deformed scene. We provide the figure with higher resolution in the supplementary.

The quantitative comparison is summarized in Table 1. We observe that: (1) Our NeuralEditor consistently outperforms all the baselines and variants for both zero-shot inference and fine-tuning settings. (2) With the precise point cloud generated by NeuralEditor, the Naive Plotting baseline even consistently outperforms PointNeRF. (3) Our ‘w/o IST’ variant has a comparable performance to Naive Plotting with the same point cloud and features in the zero-shot inference setting, but after fine-tuning its performance is significantly higher than Naive Plotting, validating the capability of NeuralEditor in NeRF modeling.

Notably, DeformingNeRF [44] performs poorly in our benchmark with significantly lower metric values. As shown in Fig. 7, the cages provided by DeformingNeRF are too coarse, and cannot even cover the whole scene. Therefore, DeformingNeRF cannot faithfully perform the precise deformation in our benchmark, leading to poor rendering results. On the contrary, both PointNeRF and our NeuralEditor at least faithfully perform the deformation, showing that point cloud is necessary for precise shape editing.

Scene Morphing Results. The morphing results between 4 NeRF Synthetic scenes are shown in Fig. 9. The

morphing process starts from Chair, morphs to Hotdog, Lego, Mic, and at last turns back to Chair. NeuralEditor produces smooth rendering results on the point cloud diffusion-guided [24] intermediate scenes, mixing the textures of the two scenes in a reasonable way. In comparison, the rendering results produced by PointNeRF are unsatisfactory, with blurry textures, black shadows, and gloomy, non-smooth colors. These results show that our NeuralEditor can render challenging intermediate morphing scenes and achieve decent results *with only the input of moved points*.

6. Conclusion

This paper proposes NeuralEditor, a point cloud-guided NeRF model that supports general shape editing tasks by manipulating the underlying point clouds. Empirical evaluation shows NeuralEditor to produce rendering results of much higher quality than baselines in a zero-shot inference manner, further significantly improving after fast fine-tuning. NeuralEditor even supports smooth scene morphing between multiple scenes, which is difficult for prior work. We hope that our work can inspire more research on point cloud-guided NeRFs and 3D shape and scene editing tasks.

Acknowledgement. This work was supported in part by NSF Grant 2106825, NIFA Award 2020-67021-32799, the Jump ARCHES endowment, the NCSA Fellows program, the IBM-Illinois Discovery Accelerator Institute, the Illinois-Inspire Partnership, and the Amazon Research Award. This work used NVIDIA GPUs at NCSA Delta through allocation CIS220014 from the ACCESS program. We thank the authors of NeRF [27] for their help in processing Blender files of the NS dataset.

References

- [1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120(2):153–168, 2016. 5
- [2] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Learning representations and generative models for 3D point clouds. In *ICML*, 2018. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 1, 2
- [4] Jon Louis Bentley. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18(9):509–517, 1975. 2, 3, 7
- [5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. MVSNeRF: Fast generalizable radiance field reconstruction from multi-view stereo. In *ICCV*, 2021. 2
- [6] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, 2019. 2
- [7] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 6
- [8] Frank Dellaert and Lin Yen-Chen. Neural volume rendering: NeRF and beyond. *arXiv:2101.05204*, 2021. 2
- [9] Kyle Gao, Yina Gao, Hongjie He, Denning Lu, Linlin Xu, and Jonathan Li. NeRF: Neural radiance field in 3D vision, a comprehensive review. *arXiv:2101.05204*, 2021. 2
- [10] Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. NeRFReN: Neural radiance fields with reflections. In *CVPR*, 2022. 2, 15
- [11] Peter Hedman, Tobias Ritschel, George Drettakis, and Gabriel Brostow. Scalable inside-out image-based rendering. *ACM Trans. Graph.*, 35(6):231:1–231:11, 2016. 2
- [12] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. DeepMVS: Learning multi-view stereopsis. In *CVPR*, 2018. 5
- [13] Alec Jacobson, Ilya Baran, Ladislav Kavan, Jovan Popović, and Olga Sorkine. Fast automatic skinning transformations. *ACM Trans. Graph.*, 31(4), 2012. 3
- [14] Tomas Jakab, Richard Tucker, Ameesh Makadia, Jiajun Wu, Noah Snavely, and Angjoo Kanazawa. Keypointdeformer: Unsupervised 3D keypoint discovery for shape control. In *CVPR*, 2021. 3
- [15] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3D neural network for multiview stereopsis. In *ICCV*, 2017. 2
- [16] Tao Ju, Scott Schaefer, and Joe Warren. Mean value coordinates for closed triangular meshes. *ACM Trans. Graph.*, 24(3):561–566, 2005. 3
- [17] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and Temples: Benchmarking large-scale scene reconstruction. *ACM Trans. Graph.*, 36(4), 2017. 6, 12, 13
- [18] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing NeRF for editing via feature field distillation. In *NeurIPS*, 2022. 2
- [19] Marc Levoy and Pat Hanrahan. Light field rendering. In *SIGGRAPH*, 1996. 2
- [20] Fayao Liu, Chunhua Shen, Guosheng Lin, and Ian Reid. Learning depth from single monocular images using deep convolutional neural fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(10):2024–2039, 2015. 2
- [21] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 2
- [22] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Junyan Zhu, and Bryan C. Russell. Editing conditional radiance fields. In *ICCV*, 2021. 2
- [23] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3D surface construction algorithm. In *SIGGRAPH*, 1987. 3
- [24] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3D point cloud generation. In *CVPR*, 2021. 2, 7, 8, 14, 16
- [25] Bruce Merry, Patrick Marais, and James Gain. Animation space: A truly linear framework for character animation. *ACM Trans. Graph.*, 25(4):1400–1423, 2006. 3
- [26] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, 2019. 2
- [27] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 4, 6, 8, 11, 17
- [28] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3D representations without 3D supervision. In *CVPR*, 2020. 2
- [29] Jesús Nieto and Toni Susin. Cage based deformations: A survey. *Lecture Notes in Computational Vision and Biomechanics*, 7:75–99, 2013. 3
- [30] Yicong Peng, Yichao Yan, Shengqi Liu, Yuhao Cheng, Shanyan Guan, Bowen Pan, Guangtao Zhai, and Xiaokang Yang. CageNeRF: Cage-based neural radiance field for generalized 3D deformation and animation. In *NeurIPS*, 2022. 2, 3, 5, 6, 7, 15
- [31] Bui Tuong Phong. Illumination for computer generated pictures. *Commun. ACM*, 18(6):311–317, 1975. 2, 4
- [32] Yi-Ling Qiao, Alexander Gao, and Ming C. Lin. Neu-Physics: Editable neural geometry and physics from monocular videos. In *NeurIPS*, 2022. 2
- [33] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinrong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxtels: Radiance fields without neural networks. In *CVPR*, 2022. 1, 2
- [34] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Photo tourism: Exploring photo collections in 3D. *ACM Trans. Graph.*, 25(3):835–846, 2006. 2

- [35] Jiaxiang Tang, Xiaokang Chen, Jingbo Wang, and Gang Zeng. Compressible-composable NeRF via rank-residual decomposition. In *NeurIPS*, 2022. 2
- [36] Jean-Marc Thiery, Julien Tierny, and Tamy Boubekeur. Jacobians and Hessians of mean value coordinates for closed triangular meshes. *The Visual Computer*, 30(9):981–995, 2014. 3
- [37] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T. Barron, and Pratul P. Srinivasan. Ref-NeRF: Structured view-dependent appearance for neural radiance fields. In *CVPR*, 2022. 1, 2, 4, 5, 14
- [38] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2Mesh: Generating 3D mesh models from single RGB images. In *ECCV*, 2018. 2
- [39] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learning multi-view image-based rendering. In *CVPR*, 2021. 1, 2
- [40] Suttisak Wizadwongsu, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwanjanakorn. NeX: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2
- [41] Liwen Wu, Jae Yong Lee, Anand Bhattacharjee, Yu-Xiong Wang, and David Forsyth. DIVeR: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. In *CVPR*, 2022. 2, 4
- [42] Dejia Xu, Peihao Wang, Yifan Jiang, Zhiwen Fan, and Zhangyang Wang. Signal processing for implicit neural representations. In *NeurIPS*, 2022. 2
- [43] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-NeRF: Point-based neural radiance fields. In *CVPR*, 2021. 2, 3, 6, 7, 11, 12, 16, 17
- [44] Tianhan Xu and Tatsuya Harada. Deforming radiance fields with cages. In *ECCV*, 2022. 2, 3, 5, 6, 7, 8, 11, 12, 15, 16, 17
- [45] 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 *ICCV*, 2021. 2
- [46] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. MVSNet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018. 5
- [47] Wang Yifan, Noam Aigerman, Vladimir G Kim, Siddhartha Chaudhuri, and Olga Sorkine-Hornung. Neural cages for detail-preserving 3D deformations. In *CVPR*, 2020. 3
- [48] Alex Yu, Rui long Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 1, 2
- [49] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [50] Yu-Jie Yuan, Yu-Kun Lai, Tong Wu, Lin Gao, and Ligang Liu. A revisit of shape editing techniques: From the geometric to the neural viewpoint. *Journal of Computer Science and Technology*, 36(3):520–554, 2021. 3
- [51] 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 *CVPR*, 2022. 2, 3, 5, 6, 7, 15
- [52] Yuzhe Zhang, Jianmin Zheng, and Yiyu Cai. Proxy-driven free-form deformation by topology-adjustable control lattice. *Computers & Graphics*, 89:167–177, 2020. 3
- [53] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 4, 14

Supplementary Material

This document contains additional descriptions (*e.g.*, implementation details, experimental setting details, *etc.*) and extra experiments (*e.g.*, ablation study, deformation on the Tanks and Temples dataset, *etc.*).

A. Black/White Background

In the main paper, we evaluated the models and presented the results on a black background, for better contrast and clearer detail visualization. Here, we provide experimental results for three representative deformed scenes of NeRF Synthetic [27] on a white background. As shown in Table A, we have the same conclusions as the experiment on a black background: NeuralEditor *significantly and consistently* outperforms PointNeRF [43] and Naive Plotting, in both zero-shot inference and fine-tuning settings.

Interestingly, we observe that the results on a white background have consistently worse metric values than those on a black background. We find that this phenomenon is *not specific* to our task of shape deformation. In fact, even in rendering the original scene, we observe that a white background results in lower metric values, as shown in Table B for the experiment of PointNeRF 20K (PointNeRF trained for 20K epochs, following an evaluation setting in [43]) on Lego of NeRF Synthetic. The impact of the background color on NeRF rendering is an interesting aspect for future investigation.

B. Additional Ablation Study

Our additional ablation study results are in Table C. We observe that:

- All components in our NeuralEditor, including infinitesimal surface transformation (IST), deterministic integration, and Phong reflection, benefit the rendering results. Notably, our NeuralEditor still outperforms PointNeRF *without* any of these components.
- For the variant without Phong reflection modeling, the zero-shot performance is close to that of the full NeuralEditor, but the performance gap becomes much larger after fine-tuning. This demonstrates that the use of visual attributes modeled by Phong reflection (*e.g.*, tint) helps fast fitting to the ambient environment.
- Our NeuralEditor’s performance drops with the initial point cloud or the final point cloud produced by PointNeRF, showing the importance of a precise point cloud optimized by our NeuralEditor for the rendering task. In these experiments, we apply de-noising on the point clouds to support IST, which is not compatible with the noisy original point clouds generated by PointNeRF.

Model	Zero-Shot Inference			Fine-Tune for 1 Epoch		
	Hotdog	Lego	Drums	Hotdog	Lego	Drums
PSNR \uparrow						
DeformingNeRF [44]	-	12.28	-	-	-	-
PointNeRF [43]	25.48	23.84	20.52	35.71	30.10	26.71
Naive Plotting	26.87	24.70	20.92	36.00	30.94	27.45
NeuralEditor w/o IST	26.88	24.71	20.92	36.27	31.45	27.64
NeuralEditor (Ours)	27.27	26.13	21.41	36.66	31.70	27.68
SSIM \uparrow						
DeformingNeRF	-	0.690	-	-	-	-
PointNeRF	0.948	0.932	0.902	0.987	0.976	0.955
Naive Plotting	0.951	0.932	0.913	0.987	0.979	0.963
NeuralEditor w/o IST	0.951	0.932	0.913	0.987	0.981	0.963
NeuralEditor (Ours)	0.957	0.964	0.920	0.988	0.983	0.964
LPIPS AlexNet \downarrow						
DeformingNeRF	-	0.271	-	-	-	-
PointNeRF	0.072	0.064	0.107	0.033	0.025	0.067
Naive Plotting	0.066	0.055	0.086	0.022	0.019	0.044
NeuralEditor w/o IST	0.064	0.054	0.085	0.021	0.017	0.043
NeuralEditor (Ours)	0.059	0.031	0.079	0.021	0.016	0.043
LPIPS VGG \downarrow						
DeformingNeRF	-	0.291	-	-	-	-
PointNeRF	0.079	0.088	0.108	0.053	0.054	0.080
Naive Plotting	0.080	0.089	0.093	0.047	0.049	0.062
NeuralEditor w/o IST	0.079	0.087	0.092	0.045	0.043	0.061
NeuralEditor (Ours)	0.073	0.056	0.087	0.044	0.041	0.060

Table A. Consistent with the results on a *black* background in the main paper, we have the same conclusions when using on a *white* background here: NeuralEditor *significantly and consistently* outperforms PointNeRF [43] and Naive Plotting on the three representative deformed scenes of NeRF Synthetic [27], in both zero-shot inference and fine-tuning settings. With the precise point cloud generated by NeuralEditor, even Naive Plotting consistently outperforms PointNeRF. The metrics investigated here are peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM), and learned perceptual image patch similarity (LPIPS).

Background Color	PSNR \uparrow
White [43]	32.40
Black	32.99

Table B. For the conventional rendering [27] of the original scene of NeRF Synthetic (*e.g.*, Lego here), the metric values are also slightly worse when using a white background as in prior work [27, 43] than a black background.

- With half of the points, the performance of NeuralEditor decreases, but it still outperforms the baseline PointNeRF. This validates that it is not the model capacity but our model design together with the precise point clouds that leads to our superior performance.
- Our NeuralEditor even supports point cloud optimiza-

Type	Variant	PSNR \uparrow on Hotdog	
		Zero-Shot	Fine-Tune
NeRF Model	Full NeuralEditor – our improved point cloud-guided NeRF + PointNeRF [43]	27.49	37.22
		25.95	36.08
Component	Full NeuralEditor – IST ('Ours w/o IST') – integration + traditional point sampling – deterministic integration + stochastic integration – NeRF modeling + plotting ('Naive Plotting') – normal vectors – Phong reflection color modeling + traditional color modeling	27.49	37.22
		27.02	36.69
		27.48	36.69
		27.46	36.87
		27.01	36.38
		27.26	37.00
		27.21	36.51
Point Cloud	Full NeuralEditor – point cloud optimization (w/ initial point cloud) – our optimized point cloud + point cloud optimized by PointNeRF – 50% points	27.49	37.22
		25.56	35.93
		26.61	35.68
		26.84	36.59
Fine-Tune	Fine-tune 0 epoch ('zero-shot')	27.49	
	Fine-tune 1 epoch	37.22	
	Fine-tune 4 epochs	38.12	
	Fine-tune 10 epochs	38.56	
	Fine-tune 10 epochs w/ point cloud optimization	38.87	

Table C. Ablation study experiments show that (1) All components in NeuralEditor benefit the rendering results on deformed scenes; (2) Our NeuralEditor generates precise point clouds, which are crucial for shape editing tasks; (3) Our NeuralEditor even supports point cloud optimization during fine-tuning to further improve the rendering performance. ‘–’ denotes that a certain component is removed, while ‘+’ denotes that a certain component is added.

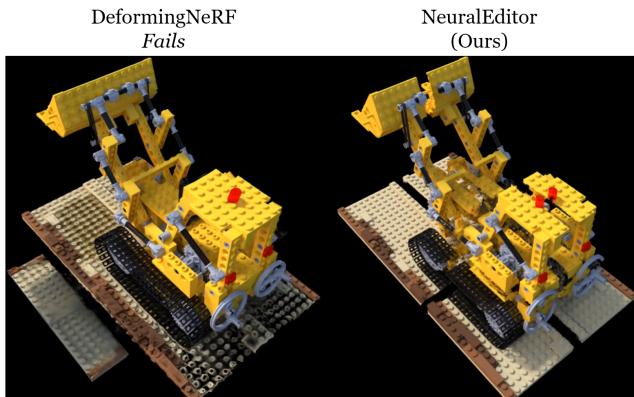


Figure A. Unlike Deforming-NeRF [44], our NeuralEditor has native support for non-continuous deformation tasks.

tion on the deformed scene, which further improves the rendering results and achieves better PSNR than fine-tuning without point cloud optimization for the same epochs. The detailed settings of fine-tuning are described in Section H.4.

C. Non-Continuous Deformation Task

Figure A shows a non-continuous deformation task con-

structed on the Lego scene of NeRF Synthetic, by cutting the scene from the xOy , yOz , and zOx planes. NeuralEditor natively supports such deformation, while DeformingNeRF fails, further validating the superiority of NeuralEditor over cage-based methods for tackling deformation tasks.

D. Experiment on Tanks and Temples Dataset

We also present the evaluation results of our NeuralEditor and baselines Naive Plotting and PointNeRF [43] on the Tanks and Temples dataset [17], as shown in Figure B. As this dataset provides white-background ground truth images for original scenes, we use a white background for evaluation and visualization. Our NeuralEditor still produces better rendering results with fewer artifacts than baselines.

Note that Tanks and Temples is not a standard NeRF dataset but a multiview stereo (MVS) dataset. It contains some background regions that are not fully cut out, e.g., the blue bar on the top of Caterpillar’s cab (2nd column in Figure B) is a part of the sky background, making the NeRF rendering results blurry and noisy as some points are mistakenly grown in those regions. We propose a background sphere technique (explained in Section I.1) to improve the rendering results in this situation. While this approach is helpful, it cannot completely resolve the issue. Additionally, inconsistent exposure settings used across dif-



Figure B. Our NeuralEditor also generates better rendering results on deformed scenes of the Tanks and Temples [17] dataset with much fewer artifacts.

ferent views of the same scene cause the rendering results to appear blurry and flashing. Therefore, all three methods cannot achieve rendering results as clean and realistic as those on NeRF Synthetic, but ours still significantly outperforms the baselines.

E. Intermediate Point Clouds for Scene Morphing

We show the intermediate point clouds for the scene morphing task (Figure 9) in Figure D.

F. Other Metrics in Table 1

We provide the results of the shape deformation task (Table 1 in the main paper) under the metrics of peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM), and learned perceptual image patch similarity (LPIPS) in Table D.

G. High-Resolution Visualization Figures

Here we provide the higher-resolution versions of the same experimental visualization figures in the main paper. The correspondence is listed below:

- Figure 6 (optimized point clouds): Figure E.
- Figure 7 (baseline DeformingNeRF): Figure F.

- Figure 8 (shape deformation): Figure G.
- Figure 9 (scene morphing): Figure H, with full morphing results for the baseline PointNeRF.

H. Implementation Details

H.1. K-D Voxels & Integration

We use a 21-layer K-D tree to build K-D voxels, which can contain up to 2^{21} (~ 2 million) points. Compared with PointNeRF that typically deals with 1 million points, our NeuralEditor can hold twice the number of points for higher capacity, and also achieves better results with the same magnitude of points as PointNeRF, as shown in Section B.

We use the voxels in the bottom 4 non-leaf layers for rendering, and use the additional 5th layer from the bottom only for point cloud growing. When doing integration for rendering, we uniformly select s points on the intersect interval, including the two side points, and apply spline integration using the features of these points, where we choose $s = 8, 11, 16, 22, 22$ for the 5 used layers from the bottom. For the feature of each selected point, we use its K nearest neighbors (KNN) for interpolation, where $K = 32$.

Such integration is not only for aggregating the average point features over the interval, but also for calculating the average normal vectors and average IST – we use KNN to interpolate these values in a similar manner as interpolating

the features. For each segment, we use the average point feature, average normal vector, and average IST in the color modeling of the representative point.

H.2. Phong Reflection Color Modeling

Consistent with RefNeRF [37], we use several multilayer perceptrons (MLPs) that directly take the integrated average feature as input to obtain the tint, roughness, modeled normal vector, and diffuse and specular color. The modeled normal vectors are trained with and controlled by the regularization losses introduced in RefNeRF. These normal vectors may be different from the normal vectors estimated from the point cloud, and they are optimized towards the estimated normal vectors via an additional regularization loss weighted by point confidence. The effect of this somewhat redundant modeling of normal vectors is two-fold: (1) The outlier points with abnormal normal vectors can be difficult to fit into the modeled ones, so their confidence values will be driven to zero when minimizing the regularization loss; (2) The estimated normal vectors can supervise and provide extra shape information for the point features through the MLP that models the normal vectors.

H.3. Point Cloud Optimization

Pruning & Growing. Following PointNeRF, we define the defects of a point cloud as two types: *outliers* and *holes*. So any shape defection condition can be decomposed into a sequence of these two types of simple defects. We design our optimization strategy based on the “pruning and growing” (P&G) method introduced in PointNeRF, which prunes the outliers by driving their confidence values to zero during training, and grows the point cloud to probe holes by selecting points from the sampled points on the training rays.

Our strategy is different from PointNeRF in several important ways. In our growing process, we grow all rays in the training dataset at a special evaluation epoch (“growing epoch”). For each ray, we pick the point far from any point in the point cloud with the highest volume density as a candidate. The candidate point whose ray has a higher pixel rendering loss is assigned a higher priority for being added to the point cloud. We then introduce a K-D tree-guided algorithm to down-sample these growing candidates: We start from the root node, and recurse on the sub-nodes until we reach a pre-set voxel size; we then add the candidate point with the highest priority in this voxel to the point cloud, while disregarding other candidate points in the same voxel. Doing so ensures the added candidate points uniformly distributed throughout the space, by preserving only the points with the highest priorities within their respective regions. We finalize the growth of candidate points by setting their features to interpolated features and their confidence values to 0.5. Note that PointNeRF uses P&G in conjunction with a stochastic training process, whereas we apply it to

our deterministic NeuralEditor during a standalone growing epoch, which enhances training stability. As a result, we obtain precise point clouds while PointNeRF cannot.

Point Cloud Denoising. It is common that a point cloud contains some noisy or isolated points. To remove these isolated points, we use Open3D [53] to identify statistical outliers w.r.t. the K -th nearest neighbor distance of each point, where $K = 64$ in our setting. Also, as mentioned above, noisy points may have irregular estimated normal vectors, and their confidence will be driven to zero with the regularization losses introduced in Section H.2.

Point Cloud Optimization Process. We regard one “point cloud optimization process” as an additional process before one training epoch, which includes (1) applying one growing epoch to obtain grown points, (2) pruning the isolated points and the points with confidence values lower than 0.1, and (3) constructing K-D voxels with the adjusted point cloud for further training. During training (Section H.4), we apply this process with a few training epochs to optimize the point cloud.

H.4. Training Settings

Training on Original Scene. We pre-train NeuralEditor on the original scene using per-scene optimization, and obtain a model that includes a precise point cloud and its points’ features for shape editing tasks. During per-scene pre-training, we start with the initial point cloud generated by the point generation network. We first train our model for 3 epochs as warm-up. From the 4th to the 12th epoch, we include an additional point cloud optimization process (Section H.3) before each training epoch. By the end of the 12th epoch, the point cloud is determined and precise enough. We keep tuning the model parameters on this underlying point cloud for up to 100 epochs, controlled with early stopping. Notably, the precise point cloud can be determined and obtained within the first 12 training epochs, even though the whole training process may take a long time.

Fine-Tuning on Deformed Scene. We apply the same training process during fine-tuning on deformed scenes. For the setting “Fine-tune 10 epochs w/ point cloud optimization” in Table C, we first train one epoch as warm-up, then train with an additional point cloud optimization process for 5 epochs, and continue to fine-tune on the optimized point cloud for the rest 4 epochs.

H.5. Matching Algorithm for Scene Morphing

To perform scene morphing, we generate the intermediate point clouds using a point cloud diffusion model [24]. However, these point clouds are *unindexed*. To render the scene with NeuralEditor, we need to assign an index to each point. This index assignment can be solved by a matching algorithm. Specifically, given point clouds P_0, P_1, \dots, P_n ,

we can match the points in each adjacent pair of point clouds P_i and P_{i+1} , and then permute the points in P_{i+1} according to the matching to align the indices.

We design a simple matching algorithm based on K-D trees. We simultaneously build two K-D trees, one for each point cloud. At each node, we select the same division axis for the two point clouds, according to their union point set. At each leaf node, we match a pair of single points from each point cloud. Our algorithm aims to match points in the two point clouds that have similar relative locations. In an ideal case where there are numerous intermediate point clouds and every adjacent pair of point clouds is sufficiently close, this algorithm will lead to highly accurate matching results.

I. Limitations

I.1. Point Cloud-Guided NeRF

Despite offering several advantages over traditional NeRFs, the current point cloud-guided NeRF models, including ours and PointNeRF, still have some limitations. First, the scene is modeled with an explicit representation (the point cloud), which is not robust when modeling surfaces with complicated visual effects, *e.g.*, a semi-transparent blurry mirror. When the model fails to interpret such visual effects, our point cloud optimization might not generate correct points close to the surface, limiting the model’s ability to improve results. Additionally, these models cannot well support strategies in NeRFReN [10] for simultaneously modeling the real-world scene and the mirrored scene to better handle mirror reflections. This is because a point cloud-guided NeRF model relies on MVS-based initialization and point cloud optimization that cannot accommodate such a form of co-optimization, which may significantly change the shape of the mirrored scene during training.

Another limitation of the point cloud-guided NeRF models is their non-robustness against inaccurate background masks, as discussed in Section D. In the existing datasets, a mask is often given to distinguish the foreground from the background, so we only need to train NeRF on the foreground. However, if the mask is not precise, some regions of the background would be mistakenly included in the mask, as in the case of the Tanks and Temples dataset (*e.g.*, the sky in the Caterpillar scene in Figure C). Since a point cloud-guided NeRF model requires points to represent the entire scene including the background, it will try to grow the points in those background regions and use them to represent the background. On the other hand, the background regions can be far from the foreground, so the model will grow the points near the scene object rather than their true location. In normal situations, each point only represents a specific region of the scene. By contrast, for each of

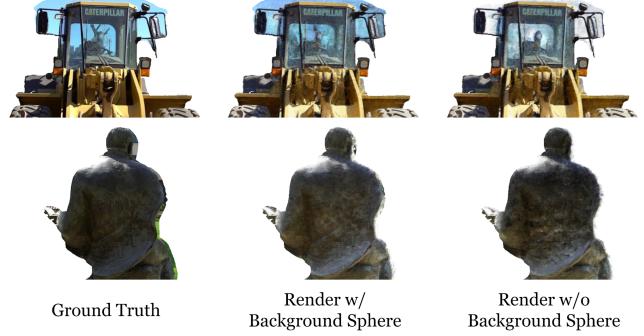


Figure C. Modeling a background sphere helps NeuralEditor to differentiate between the background and the foreground scene, thus preventing it from growing wrong points to model the background. Doing so potentially improves the robustness of the point cloud-guided NeRF model.

these mistakenly grown points, its view-dependent color is trained with different regions of the background which are inconsistent, thus resulting in abnormal rendering results in novel views.

To address this issue, we enhanced NeuralEditor by modeling a *background sphere* that covers the entire scene. Therefore, the point cloud-guided NeRF model can directly represent the scene’s background using the sphere, instead of growing new points. Such a strategy was introduced for the Tanks and Temples dataset in Figure B. Figure C further analyzes the impact of this background sphere, validating its effectiveness in approximately modeling background regions without growing wrong points. Further investigation is worthwhile in other strategies to tackle this issue.

I.2. Environment Modeling in Shape Deformation Task

Neither our work nor existing methods [30, 44, 51] take into account the surrounding ambient environment when addressing the shape deformation task. These methods thus cannot assign different colors to the scene according to the changes in lighting conditions, as shown in the top of the shovel and the bent chimney in the Lego scene and the shadow on the cushion in the Chair scene (Figure G). Fortunately, our NeuralEditor, incorporating the Phong reflection’s visual attributes (*e.g.*, tint), enables fast fitting to the ambient environment through fine-tuning on the deformed scene (Section B).

I.3. Evaluation for Scene Morphing Task

Quantitatively evaluating the visual realism of intermediate scenes during morphing is challenging, due to the lack of ground truth and associated metrics, as these scenes “do not exist” in the real world. Therefore, we rely mainly on visualizations for evaluation.

Model	Zero-Shot Inference								Fine-Tune for 1 Epoch							
	Chair	Hotdog	Lego	Drums	Ficus	Materials	Mic	Ship	Chair	Hotdog	Lego	Drums	Ficus	Materials	Mic	Ship
PSNR \uparrow																
DeformingNeRF [44]	18.84	-	13.10	-	-	-	-	-	-	-	-	-	-	-	-	-
PointNeRF [43]	22.21	25.95	24.56	21.00	24.24	21.21	26.77	21.19	30.11	36.08	31.45	27.16	31.48	27.55	34.34	28.90
Naive Plotting	24.91	27.01	25.64	21.29	26.22	21.65	27.63	22.29	32.01	36.38	31.72	28.09	33.21	30.31	35.15	30.01
NeuralEditor w/o IST	24.92	27.02	25.65	21.29	26.24	21.64	27.64	22.28	32.24	36.69	32.79	28.30	33.34	30.40	35.28	30.08
NeuralEditor (Ours)	25.85	27.49	27.46	21.84	27.19	23.18	27.75	24.16	32.53	37.22	32.95	28.35	33.53	30.82	35.46	30.44
SSIM \uparrow																
DeformingNeRF	0.865	-	0.645	-	-	-	-	-	-	-	-	-	-	-	-	-
PointNeRF	0.910	0.947	0.933	0.890	0.915	0.856	0.945	0.759	0.972	0.988	0.977	0.953	0.973	0.925	0.981	0.875
Naive Plotting	0.950	0.954	0.934	0.908	0.953	0.887	0.964	0.844	0.984	0.987	0.977	0.962	0.985	0.968	0.988	0.935
NeuralEditor w/o IST	0.950	0.954	0.934	0.908	0.953	0.887	0.964	0.845	0.984	0.988	0.982	0.963	0.985	0.968	0.988	0.936
NeuralEditor (Ours)	0.963	0.960	0.966	0.916	0.960	0.909	0.966	0.875	0.986	0.989	0.983	0.963	0.986	0.971	0.989	0.939
LPIPS AlexNet \downarrow																
DeformingNeRF	0.071	-	0.297	-	-	-	-	-	-	-	-	-	-	-	-	-
PointNeRF	0.049	0.080	0.067	0.122	0.062	0.099	0.057	0.139	0.018	0.033	0.023	0.077	0.028	0.067	0.035	0.073
Naive Plotting	0.038	0.063	0.053	0.088	0.047	0.098	0.039	0.159	0.014	0.023	0.019	0.047	0.021	0.042	0.020	0.081
NeuralEditor w/o IST	0.037	0.062	0.052	0.087	0.046	0.097	0.039	0.158	0.013	0.021	0.015	0.046	0.020	0.041	0.019	0.080
NeuralEditor (Ours)	0.030	0.057	0.029	0.080	0.042	0.076	0.037	0.126	0.012	0.020	0.015	0.045	0.019	0.035	0.019	0.075
LPIPS VGG \downarrow																
DeformingNeRF	0.067	-	0.291	-	-	-	-	-	-	-	-	-	-	-	-	-
PointNeRF	0.047	0.082	0.060	0.115	0.070	0.091	0.044	0.153	0.019	0.055	0.055	0.086	0.039	0.061	0.029	0.090
Naive Plotting	0.051	0.080	0.091	0.094	0.070	0.109	0.040	0.183	0.029	0.050	0.049	0.068	0.044	0.065	0.030	0.129
NeuralEditor w/o IST	0.051	0.079	0.088	0.093	0.069	0.107	0.040	0.182	0.027	0.048	0.043	0.066	0.042	0.064	0.029	0.127
NeuralEditor (Ours)	0.041	0.074	0.057	0.088	0.067	0.095	0.038	0.163	0.026	0.045	0.042	0.065	0.042	0.058	0.028	0.121

Table D. **Full comparison results of Table 1** in the main paper under all metrics. NeuralEditor significantly and consistently outperforms PointNeRF and Naive Plotting on all deformed scenes of NeRF Synthetic *under all metrics*, in both zero-shot inference and fine-tuning settings. Our infinitesimal surface transformation (IST) effectively improves the results by correcting the view-dependent colors. With the precise point cloud generated by NeuralEditor, even Naive Plotting consistently outperforms PointNeRF.

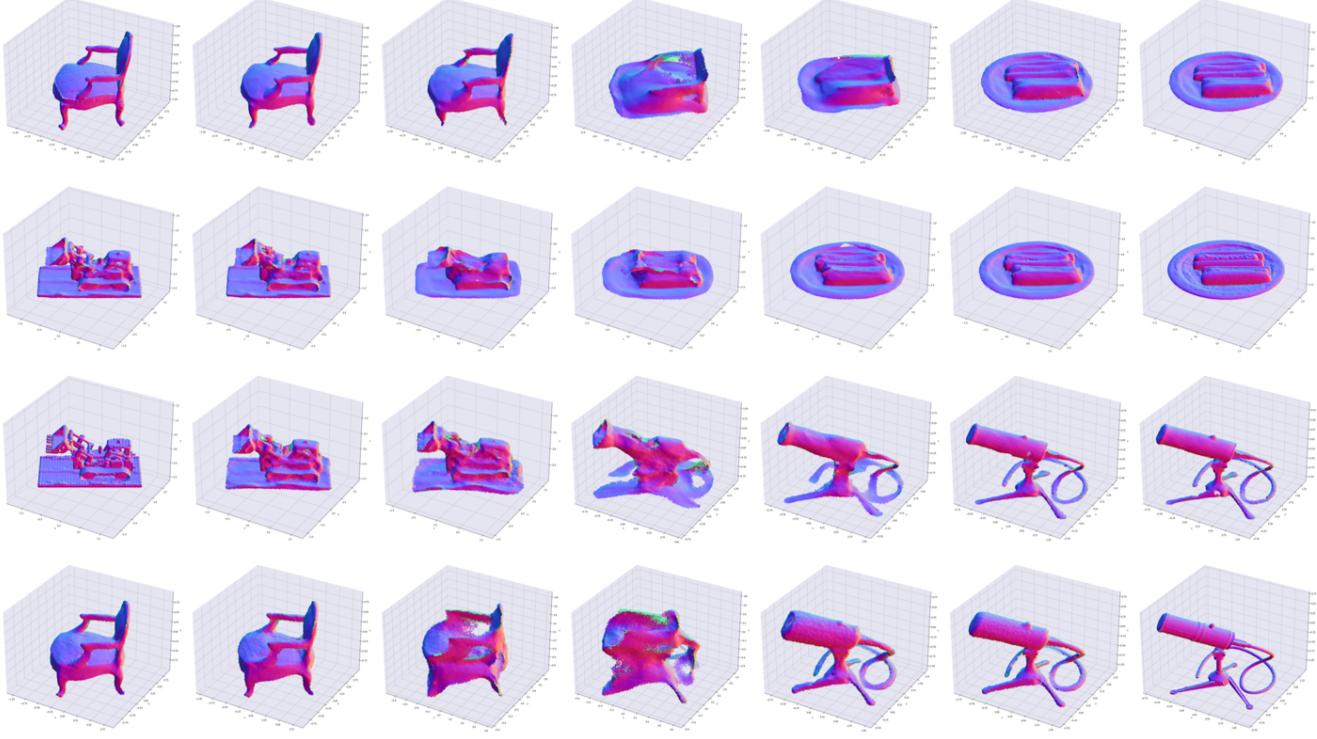


Figure D. Intermediate point clouds for the scene morphing task corresponding to Figure 9 in the main paper, which are generated by the point cloud diffusion model [24].

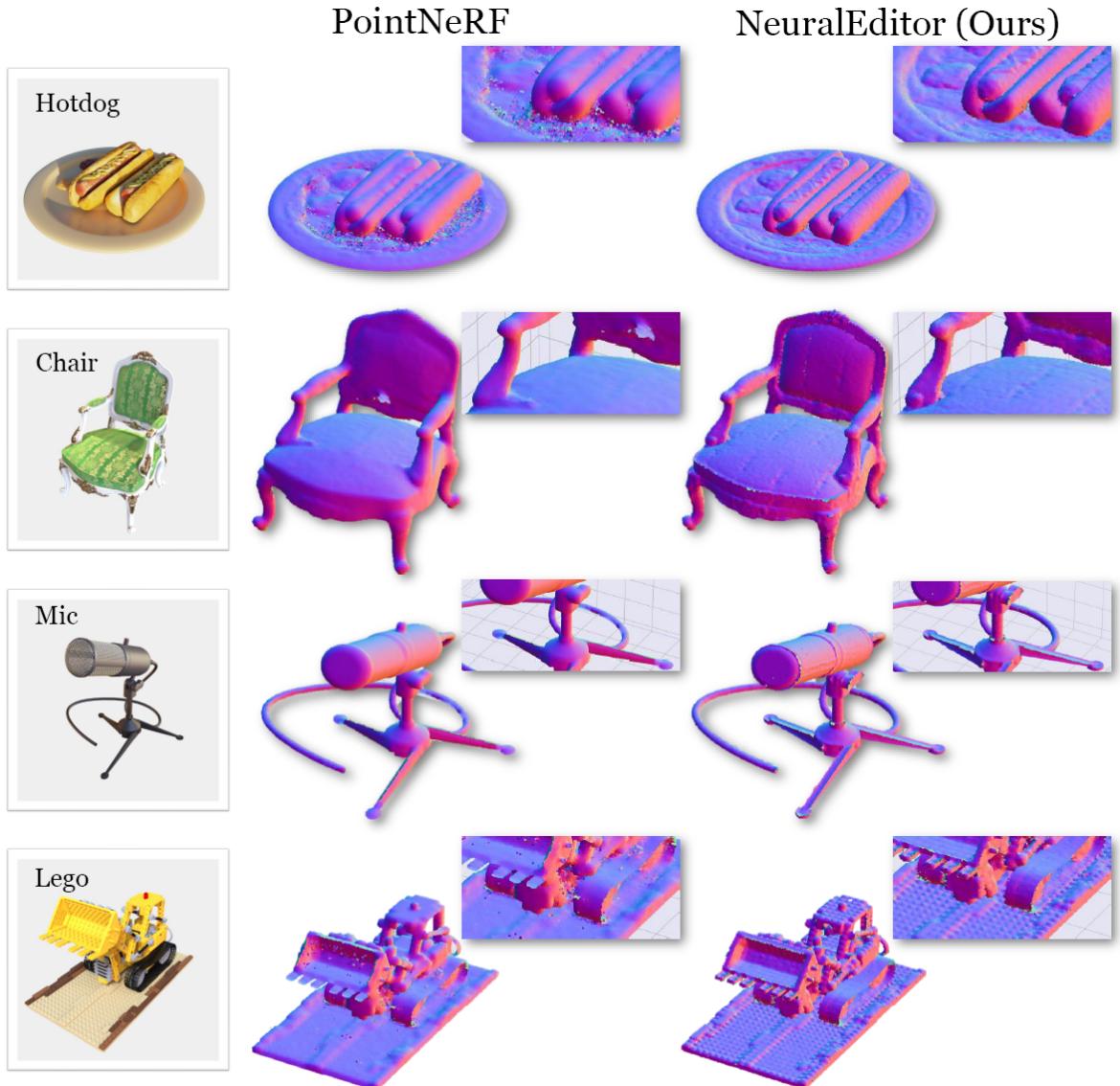


Figure E. **High-resolution version of Figure 6** in the main paper for more detailed visualization. NeuralEditor generates much more precise point clouds than PointNeRF [43] in the four scenes of NeRF Synthetic [27]. The points are colored with their normal vectors.

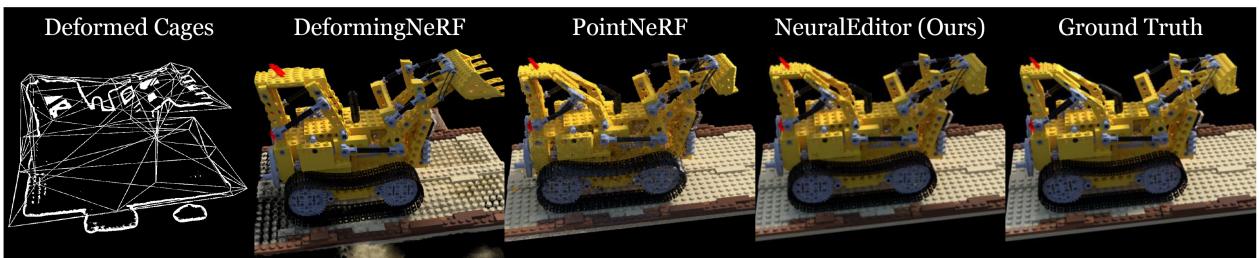


Figure F. **High-resolution version of Figure 7** in the main paper for more detailed visualization. With too coarse cages, DeformingNeRF [44] is unable to perform the deformation faithfully, leading to poor results.

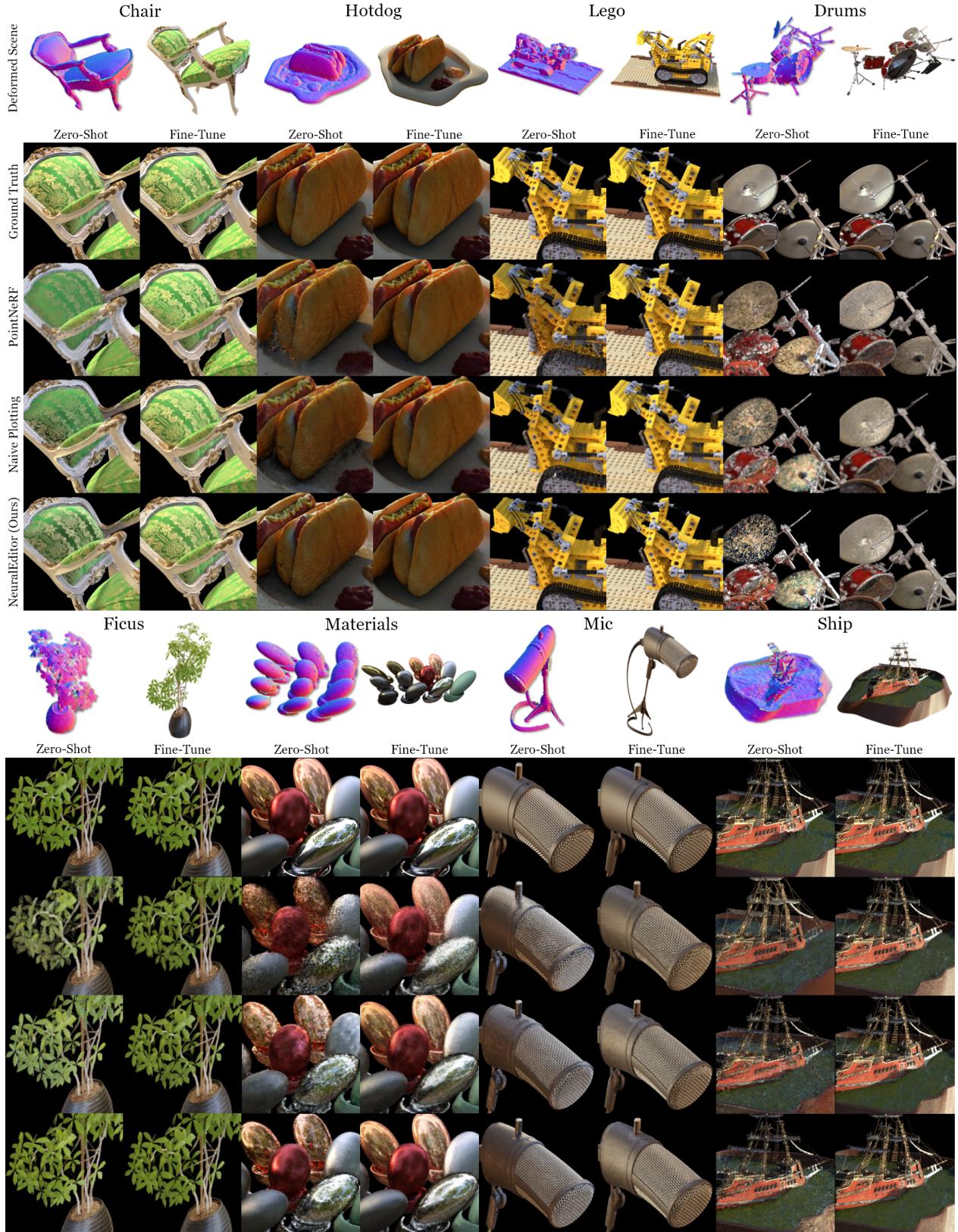


Figure G. **High-resolution version of Figure 8** in the main paper for more detailed visualization. NeuralEditor produces superior rendering results to PointNeRF, with significantly fewer artifacts in zero-shot inference. Fine-tuning further improves the *consistency of rendering with the ambient environment*. We use a black background for better visualization.



Figure H. **High-resolution, extended version of Figure 9** in the main paper for more detailed visualization. Our NeuralEditor produces *smooth morphing* results between Chair, Hotdog, Lego, and Mic in the NeRF Synthetic dataset, while PointNeRF produces results with blurry textures, black shadows, and gloomy, non-smooth colors. The rendering results in the looped morphing process are arranged in the shape of the numerical digit “3,” indicated by the dividing lines and arrows. For the baseline PointNeRF, in the main paper we showed the morphing results between Chair and Hotdog due to limited space; here we include the full morphing results across all 4 scenes.