# Direct Voxel Grid Optimization:
# Super-fast Convergence for Radiance Fields Reconstruction

Cheng Sun        Min Sun        Hwann-Tzong Chen
National Tsing Hua University

chengsun@gapp.nthu.edu.tw        sunmin@ee.nthu.edu.tw        htchen@cs.nthu.edu.tw

## Abstract

*We present a super-fast convergence approach to reconstructing the per-scene radiance field from a set of images that capture the scene with known poses. This task, which is often applied to novel view synthesis, is recently revolutionized by Neural Radiance Field (NeRF) for its state-of-the-art quality and flexibility. However, NeRF and its variants require a lengthy training time ranging from hours to days for a single scene. In contrast, our approach achieves NeRF-comparable quality and converges rapidly from scratch in less than 15 minutes with a single GPU. We adopt a representation consisting of a density voxel grid for scene geometry and a feature voxel grid with a shallow network for complex view-dependent appearance. Modeling with explicit and discretized volume representations is not new, but we propose two simple yet non-trivial techniques that contribute to fast convergence speed and high-quality output. First, we introduce the post-activation interpolation on voxel density, which is capable of producing sharp surfaces in lower grid resolution. Second, direct voxel density optimization is prone to suboptimal geometry solutions, so we robustify the optimization process by imposing several priors. Finally, evaluation on five inward-facing benchmarks shows that our method matches, if not surpasses, NeRF's quality, yet it only takes about 15 minutes to train from scratch for a new scene. Code: https://github.com/sunset1995/DirectVoxGO.*
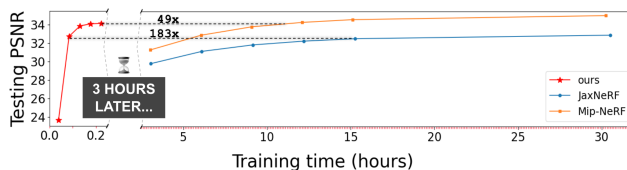
## 1. Introduction

Achieving free-viewpoint navigation of 3D objects or scenes from only a set of calibrated images as input is a demanding task. For instance, it enables online product showcase to provide an immersive user experience comparing to static image demonstration. Recently, Neural Radiance Fields (NeRFs) [37] have emerged as powerful representations yielding state-of-the-art quality on this task.

Despite its effectiveness in representing scenes, NeRF is known to be hampered by the need of lengthy training time



23.64 PSNR.        32.72 PSNR.        34.22 PSNR.
Ours at 2.33 mins.        Ours at 5.07 mins.        Ours at 13.72 mins.

(a) The synthesized novel view by our method at three training checkpoints.



(b) The training curves of different methods on *Lego* scene. The training time of each method is measured on our machine with a single NVIDIA RTX 2080 Ti GPU.

Figure 1. **Super-fast convergence by our method.** The key to our speedup is to optimize the volume density modeled in a dense voxel grid directly. Note that our method needs neither a conversion step from any trained implicit model (*e.g.*, NeRF) nor a cross-scene pretraining, *i.e.*, our voxel grid representation is directly and efficiently trained from scratch for each scene.

and the inefficiency in rendering new views. This makes NeRF infeasible for many application scenarios. Several follow-up methods [15, 18, 29, 30, 43, 44, 67] have shown significant speedup of FPS in testing phase, some of which even achieve real-time rendering. However, only few methods show training times speedup, and the improvements are not comparable to ours [1, 10, 31] or lead to worse quality [6, 60]. On a single GPU machine, several hours of per scene optimization or a day of pretraining is typically required.

To reconstruct a volumetric scene representation from a set of images, NeRF uses multilayer perceptron (MLP) to implicitly learn the mapping from a queried 3D point (with a viewing direction) to its colors and densities. The queried properties along a camera ray can then be accumulated into a pixel color by volume rendering techniques. Our work takes inspiration from the recent success [15, 18, 67] that uses classic voxel grid to explicitly store the scene properties,

1

which enables real-time rendering and shows good quality. However, their methods can not train from scratch and need a conversion step from the trained implicit model, which causes a bottleneck to the training time.

The key to our speedup is to use a dense voxel grid to directly model the 3D geometry (volume density). Developing an elaborate strategy for view-dependent colors is not in the main scope of this paper, and we simply use a hybrid representation (feature grid with shallow MLP) for colors.

Directly optimizing the density voxel grid leads to superfast converges but is prone to suboptimal solutions, where our method allocates "cloud" at free space and tries to fit the photometric loss with the cloud instead of searching a geometry with better multi-view consistency. Our solution to this problem is simple and effective. First, we initialize the density voxel grid to yield opacities very close to zero everywhere to avoid the geometry solutions being biased toward the cameras' near planes. Second, we give a lower learning rate to voxels visible to fewer views, which can avoid redundant voxels that are allocated just for explaining the observations from a small number of views. We show that the proposed solutions can successfully avoid the suboptimal geometry and work well on the five datasets.

Using the voxel grid to model volume density still faces a challenge in scalability. For parsimony, our approach automatically finds a BBox tightly encloses the volume of interest to allocate the voxel grids. Besides, we propose post-activation—applying all the activation functions after trilinearly interpolating the density voxel grid. Previous work either interpolates the voxel grid for the activated opacity or uses nearest-neighbor interpolation, which results in a smooth surface in each grid cell. Conversely, we prove mathematically and empirically that the proposed post-activation can model (beyond) a sharp linear surface within a single grid cell. As a result, we can use fewer voxels to achieve better qualities—our method with $160^3$ dense voxels already outperforms NeRF in most cases.

In summary, we have two main technical contributions. First, we implement two priors to avoid suboptimal geometry in direct voxel density optimization. Second, we propose the *post-activated voxel-grid interpolation*, which enables sharp boundary modeling in lower grid resolution. The resulting key merits of this work are highlighted as follows:

- Our convergence speed is about two orders of magnitude faster than NeRF—reducing training time from $10-20$ hours to 15 minutes on our machine with a single NVIDIA RTX 2080 Ti GPU, as shown in Fig. 1.

- We achieve visual quality comparable to NeRF at a rendering speed that is about $45\times$ faster.

- Our method does not need cross-scene pretraining.

- Our grid resolution is about $160^3$, while the grid resolution in previous work [15, 18, 67] ranges from $512^3$

to $1300^3$ to achieve NeRF-comparable quality.

## 2. Related work

**Representations for novel view synthesis.** Images synthesis from novel viewpoints given a set of images capturing the scene is a long-standing task with rich studies. Previous work has presented several scene representations reconstructed from the input images to synthesize the unobserved viewpoints. Lumigraph [4, 16] and light field representation [7, 23, 24, 47] directly synthesize novel views by interpolating the input images but require very dense scene capture. Layered depth images [11, 46, 48, 58] work for sparse input views but rely on depth maps or estimated depth with sacrificed quality. Mesh-based representations [8, 55, 59, 64] can run in real-time but have a hard time with gradient-based optimization without template meshes provided. Recent approaches employ 2D/3D Convolutional Neural Network (CNNs) to estimate multiplane images (MPIs) [12, 26, 36, 52, 57, 72] for forward-facing captures; estimate voxel grid [17, 32, 49] for inward-facing captures. Our method uses gradient-descent to optimize voxel grids directly and does not rely on neural networks to predict the grid values, and we still outperform the previous works [17, 32, 49] with CNNs by a large margin.

**Neural radiance fields.** Recently, NeRF [37] stands out to be a prevalent method for novel view synthesis with rapid progress, which takes a moderate number of input images with known camera poses. Unlike traditional explicit and discretized volumetric representations (*e.g.*, voxel grids and MPIs), NeRF uses coordinate-based multilayer perceptrons (MLP) as an implicit and continuous volumetric representation. NeRF achieves appealing quality and has good flexibility with many follow-up extensions to various setups, *e.g.*, relighting [2, 3, 51, 71], deformation [13, 38, 40, 41, 56], self-calibration [19, 27, 28, 35, 62], meta-learning [53], dynamic scene modeling [14, 25, 33, 42, 65], and generative modeling [5, 22, 45]. Nevertheless, NeRF has unfavorable limitations of lengthy training progress and slow rendering speed. In this work, we mainly follow NeRF's original setup, while our method can optimize the volume density explicitly encoded in a voxel grid to speed up both training and testing by a large margin with comparable quality.

**Hybrid volumetric representations.** To combine NeRF's implicit representation and traditional grid representations, the coordinate-based MLP is extended to also conditioning on the local feature in the grid. Recently, hybrid voxels [18, 30] and MPIs [63] representations have shown success in fast rendering speed and result quality. We use hybrid representation to model view-dependent color as well.

**Fast NeRF rendering.** NSVF [30] uses octree in its hybrid representation to avoid redundant MLP queries in free
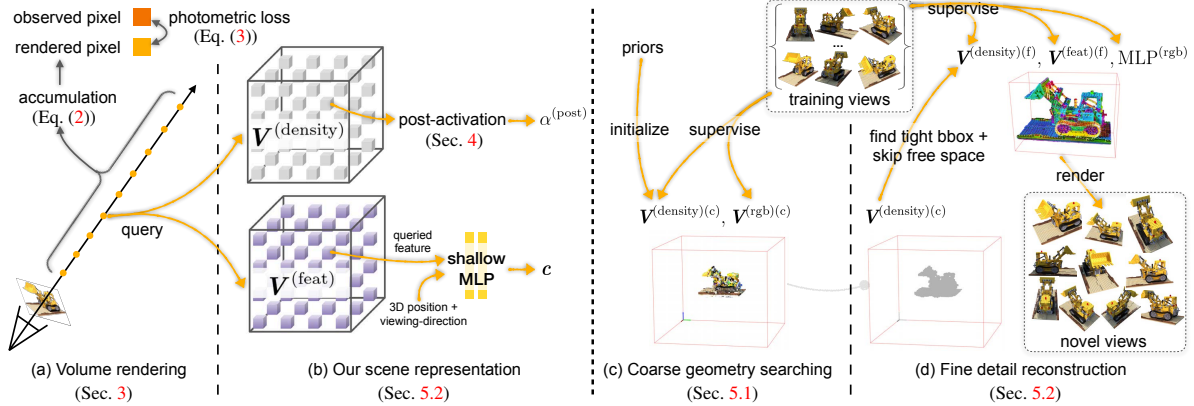
Figure 2. **Approach overview.** We first review NeRF in Sec. 3. In Sec. 4, we present a novel post-activated density voxel grid to support sharp surface modeling in lower grid resolutions. In Sec. 5, we show our approach to the reconstruction of radiance field with super-fast convergence, where we first find a coarse geometry in Sec. 5.1 and then reconstruct the fine details and view-dependent effects in Sec. 5.2.

space. However, NSVF still needs many training hours due to the deep MLP in its representation. Recent methods further use thousands of tiny MLPs [44] or explicit volumetric representations [15,18,63,67] to achieve real-time rendering. Unfortunately, gradient-based optimization is not directly applicable to their methods due to their topological data structures or the lack of priors. As a result, these methods [15, 18, 44, 63, 67] still need a conversion step from a trained implicit model (*e.g.*, NeRF) to their final representation that supports real-time rendering. Their training time is still burdened by the lengthy implicit model optimization.

**Fast NeRF convergence.** Recent works that focus on fewer input views setup also bring faster convergence as a side benefit. These methods rely on generalizable pre-training [6, 60, 68] or external MVS depth information [10, 31], while ours does not. Further, they still require several per-scene fine-tuning hours [10] or fail to achieve NeRF quality in the full input-view setup [6,60,68]. Most recently, NeuRay [31] shows NeRF's quality with 40 minutes per-scene training time in the lower-resolution setup. Under the same GPU spec, our method achieves NeRF's quality in 15 minutes per scene on the high-resolution setup and does not require depth guidance and cross-scene pre-training.

## 3. Preliminaries

To represent a 3D scene for novel view synthesis, Neural Radiance Fields (NeRFs) [37] employ multilayer perceptron (MLP) networks to map a 3D position $x$ and a viewing direction $d$ to the corresponding density $\sigma$ and view-dependent color emission $c$:

$$(\sigma, e) = \mathrm{MLP}^{(\mathrm{pos})}(x) , \quad (1a)$$

$$c = \mathrm{MLP}^{(\mathrm{rgb})}(e, d) , \quad (1b)$$

where the learnable MLP parameters are omitted, and $e$ is an intermediate embedding to help the much shallower

$\mathrm{MLP}^{(\mathrm{rgb})}$ to learn $c$ (see NeRF++ [69] for more discussions on the architecture design). In practice, positional encoding is applied to $x$ and $d$, which enables the MLPs to learn the high-frequency details from low-dimensional input [54]. For output activation, Sigmoid is applied on $c$; ReLU or Softplus is applied on $\sigma$ (see Mip-NeRF [1] for more discussion on output activation).

To render the color of a pixel $\hat{C}(r)$, we cast the ray $r$ from the camera center through the pixel; $K$ points are then sampled on $r$ between the pre-defined near and far planes; the $K$ ordered sampled points are then used to query for their densities and colors $\{(\sigma_i, c_i)\}_{i=1}^{K}$ (MLPs are queried in NeRF). Finally, the $K$ queried results are accumulated into a single color with the volume rendering quadrature in accordance with the optical model given by Max [34]:

$$\hat{C}(r) = \left( \sum_{i=1}^{K} T_i \alpha_i c_i \right) + T_{K+1} c_{\mathrm{bg}} , \quad (2a)$$

$$\alpha_i = \mathrm{alpha}(\sigma_i, \delta_i) = 1 - \exp(-\sigma_i \delta_i) , \quad (2b)$$

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) , \quad (2c)$$

where $\alpha_i$ is the probability of termination at the point $i$; $T_i$ is the accumulated transmittance from the near plane to point $i$; $\delta_i$ is the distance to the adjacent sampled point, and $c_{bg}$ is a pre-defined background color.

Given the training images with known poses, NeRF model is trained by minimizing the photometric MSE between the observed pixel color $C(r)$ and the rendered color $\hat{C}(r)$:

$$\mathcal{L}_{\mathrm{photo}} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left\| \hat{C}(r) - C(r) \right\|_2^2 , \quad (3)$$

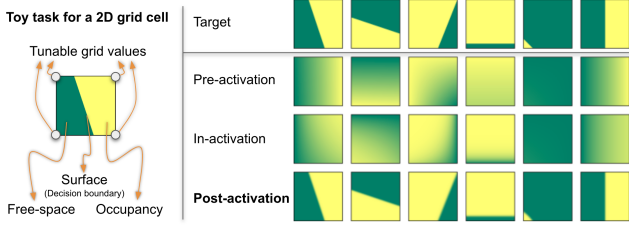where $\mathcal{R}$ is the set of rays in a sampled mini-batch.

3

Figure 3. **A single grid cell with post-activation is capable of modeling sharp linear surfaces.** *Left*: We depict the toy task for a 2D grid cell, where a grid cell is optimized for the linear surface (decision boundary) across it. *Right*: Each column shows an example task for three different methods. The results show that a single grid cell with post-activation (Eq. (6c)) is adequate to recover faithfully the linear surface. Conversely, pre-activation (Eq. (6a)) and in-activation (Eq. (6b)) fail to accomplish the tasks as they can only fit into smooth results, and thus would require more grid cells to recover the surface detail. See supplementary material for the mathematical proof.

## 4. Post-activated density voxel grid

**Voxel-grid representation.** A voxel-grid representation models the modalities of interest (*e.g.*, density, color, or feature) explicitly in its grid cells. Such an explicit scene representation is efficient to query for any 3D positions via interpolation:

$$\text{interp}(\boldsymbol{x}, \boldsymbol{V}) : \left(\mathbb{R}^3, \mathbb{R}^{C \times N_x \times N_y \times N_z}\right) \to \mathbb{R}^C , \quad (4)$$
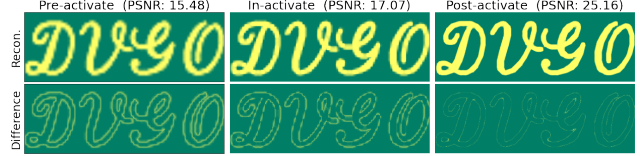
where $\boldsymbol{x}$ is the queried 3D point, $\boldsymbol{V}$ is the voxel grid, $C$ is the dimension of the modality, and $N_x \cdot N_y \cdot N_z$ is the total number of voxels. Trilinear interpolation is applied if not specified otherwise.

**Density voxel grid for volume rendering.** Density voxel grid, $\boldsymbol{V}^{(\text{density})}$, is a special case with $C = 1$, which stores the density values for volume rendering (Eq. (2)). We use $\ddot{\sigma} \in \mathbb{R}$ to denote the raw voxel density before applying the density activation (*i.e.*, a mapping of $\mathbb{R} \to \mathbb{R}_{\geq 0}$). In this work, we use the shifted softplus mentioned in Mip-NeRF [1] as the density activation:

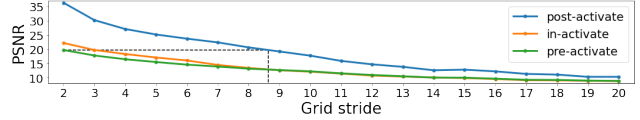$$\sigma = \text{softplus}(\ddot{\sigma}) = \log(1 + \exp(\ddot{\sigma} + b)) , \quad (5)$$

where the shift $b$ is a hyperparameter. Using softplus instead of ReLU is crucial to optimize voxel density directly, as it is irreparable when a voxel is falsely set to a negative value with ReLU as the density activation. Conversely, softplus allows us to explore density very close to 0.

**Sharp decision boundary via post-activation.** The interpolated voxel density is processed by softplus (Eq. (5)) and alpha (Eq. (2b)) functions sequentially for volume rendering. We consider three different orderings—pre-activation, in-activation, and post-activation—of plugging in the trilinear interpolation and performing the activation, given a



(a) Visual comparison of image fitting results under grid resolution $(H/5) \times (W/5)$. The first row is the results of pre-, in-, and post-activation. The second row is their per-pixel absolute difference to the target image.



(b) PSNRs achieved by pre-, in- and post-activation under different grid strides. A grid stride $s$ means that the grid resolution is $(H/s) \times (W/s)$. The black dashed line highlights that post-activation with stride $\approx 8.5$ can achieve the same PSNR as pre-activation with stride 2 in this example.

Figure 4. **Toy example on image fitting.** The target 2D image is binary to imitate the scenario that most of the 3D space is either occupied or free. The objective is to reconstruct the target image by a low-resolution 2D grid. In each optimization step, the tunable 2D grid is queried by interpolation with pre-activation (Eq. (6a)), in-activation (Eq. (6b)), or post-activation (Eq. (6c)) to minimize the mean squared error to the target image. The result reveals that the *post-activation* can produce sharp boundaries even with low grid resolution (Fig. 4a) and is much better than the other two under various grid resolutions (Fig. 4b). This motivates us to model the 3D geometry directly via voxel grids with *post-activation*.

queried 3D point $\boldsymbol{x}$:

$$\alpha^{(\text{pre})} = \underline{\text{interp}} \left(\boldsymbol{x}, \underline{\text{alpha}} \left(\underline{\text{softplus}} \left(\boldsymbol{V}^{(\text{density})}\right)\right)\right) , \quad (6a)$$

$$\alpha^{(\text{in})} = \underline{\text{alpha}} \left(\text{interp} \left(\boldsymbol{x}, \underline{\text{softplus}} \left(\boldsymbol{V}^{(\text{density})}\right)\right)\right) , \quad (6b)$$

$$\alpha^{(\text{post})} = \underline{\text{alpha}} \left(\underline{\text{softplus}} \left(\text{interp} \left(\boldsymbol{x}, \boldsymbol{V}^{(\text{density})}\right)\right)\right) . \quad (6c)$$

The input $\delta$ to the function alpha (Eq. (2b)) is omitted for simplicity. We show that the post-activation, *i.e.*, applying all the non-linear activation after the trilinear interpolation, is capable of producing sharp surfaces (decision boundaries) with much fewer grid cells. In Fig. 3, we use a 2D grid cell as an example to show that a grid cell with post-activation can produce a sharp linear boundary, while pre- and in-activation can only produce smooth results and thus require more cells for the surface detail. In Fig. 4, we further use binary image regression as a toy example to compare their capability, which also shows that post-activation can achieve a much better efficiency in grid cell usage.

## 5. Fast and direct voxel grid optimization

We depict an overview of our approach in Fig. 2. In Sec. 5.1, we first search the coarse geometry of a scene. In Sec. 5.2, we then reconstruct the fine detail including view-dependent effects. Hereinafter we use superscripts (c) and (f) to denote variables in the coarse and fine stages.

## 5.1. Coarse geometry searching

Typically, a scene is dominated by free space (*i.e.*, unoccupied space). Motivated by this fact, we aim to efficiently find the coarse 3D areas of interest before reconstructing the fine detail and view-dependent effect that require more computation resources. We can thus greatly reduce the number of queried points on each ray in the later fine stage.

**Coarse scene representation.** We use a coarse density voxel grid $V^{(\text{density})(\text{c})} \in \mathbb{R}^{1 \times N_x^{(\text{c})} \times N_y^{(\text{c})} \times N_z^{(\text{c})}}$ with post-activation (Eq. (6c)) to model scene geometry. We only model view-invariant color emissions by $V^{(\text{rgb})(\text{c})} \in \mathbb{R}^{3 \times N_x^{(\text{c})} \times N_y^{(\text{c})} \times N_z^{(\text{c})}}$ in the coarse stage. A query of any 3D point $x$ is efficient with interpolation:

$$\ddot{\sigma}^{(\text{c})} = \text{interp}\left(x, V^{(\text{density})(\text{c})}\right) , \qquad (7\text{a})$$

$$c^{(\text{c})} = \text{interp}\left(x, V^{(\text{rgb})(\text{c})}\right) , \qquad (7\text{b})$$

where $c^{(\text{c})} \in \mathbb{R}^3$ is the view-invariant color and $\ddot{\sigma}^{(\text{c})} \in \mathbb{R}$ is the raw volume density.

**Coarse voxels allocation.** We first find a bounding box (BBox) tightly enclosing the camera frustums of training views (See the red BBox in Fig. 2c for an example). Our voxel grids are aligned with the BBox. Let $L_x^{(\text{c})}, L_y^{(\text{c})}, L_z^{(\text{c})}$ be the lengths of the BBox and $M^{(\text{c})}$ be the hyperparameter for the expected total number of voxels in the coarse stage. The voxel size is $s^{(\text{c})} = \sqrt[3]{L_x^{(\text{c})} \cdot L_y^{(\text{c})} \cdot L_z^{(\text{c})}/M^{(\text{c})}}$, so there are $N_x^{(\text{c})}, N_y^{(\text{c})}, N_z^{(\text{c})} = \lfloor L_x^{(\text{c})}/s^{(\text{c})} \rfloor, \lfloor L_y^{(\text{c})}/s^{(\text{c})} \rfloor, \lfloor L_z^{(\text{c})}/s^{(\text{c})} \rfloor$ voxels on each side of the BBox.

**Coarse-stage points sampling.** On a pixel-rendering ray, we sample query points as

$$x_0 = o + t^{(\text{near})}d , \qquad (8\text{a})$$

$$x_i = x_0 + i \cdot \delta^{(\text{c})} \cdot \frac{d}{\|d\|^2} , \qquad (8\text{b})$$

where $o$ is the camera center, $d$ is the ray-casting direction, $t^{(\text{near})}$ is the camera near bound, and $\delta^{(\text{c})}$ is a hyperparameter for the step size that can be adaptively chosen according to the voxel size $s^{(c)}$. The query index $i$ ranges from 1 to $\lceil t^{(\text{far})} \cdot \|d\|^2 / \delta^{(\text{c})} \rceil$, where $t^{(\text{far})}$ is the camera far bound, so the last sampled point stops nearby the far plane.

**Prior 1: low-density initialization.** At the start of training, the importance of points far from a camera is downweighted due to the accumulated transmittance term in Eq. (2c). As a result, the coarse density voxel grid $V^{(\text{density})(\text{c})}$ could be accidentally trapped into a suboptimal "cloudy" geometry with higher densities at camera near planes. We thus have to initialize $V^{(\text{density})(\text{c})}$ more carefully to ensure that all sampled points on rays are visible to the cameras at the beginning, *i.e.*, the accumulated transmittance rates $T_i$s in Eq. (2c) are close to 1.

In practice, we initialize all grid values in $V^{(\text{density})(\text{c})}$ to 0 and set the bias term in Eq. (5) to

$$b = \log\left(\left(1 - \alpha^{(\text{init})(\text{c})}\right)^{-\frac{1}{s^{(\text{c})}}} - 1\right) , \qquad (9)$$

where $\alpha^{(\text{init})(\text{c})}$ is a hyperparameter. Thereby, the accumulated transmittance $T_i$ is decayed by $1 - \alpha^{(\text{init})(\text{c})} \approx 1$ for a ray that traces forward a distance of a voxel size $s^{(\text{c})}$. See supplementary material for the derivation and proof.

**Prior 2: view-count-based learning rate.** There could be some voxels visible to too few training views in real-world capturing, while we prefer a surface with consistency in many views instead of a surface that can only explain few views. In practice, we set different learning rates for different grid points in $V^{(\text{density})(\text{c})}$. For each grid point indexed by $j$, we count the number of training views $n_j$ to which point $j$ is visible, and then scale its base learning rate by $n_j/n_{\max}$, where $n_{\max}$ is the maximum view count over all grid points.

**Optimizing the coarse scene representation.** In addition to the photometric loss in Eq. (3), we further incorporate per-point rgb loss and background entropy loss. The per-point rgb loss is a modification of Eq. (3) that directly supervises each of the $K$ ordered sampled points instead of the accumulated ray:

$$\mathcal{L}_{\text{pt\_rgb}} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \sum_{i=1}^{K} \left(T_i \alpha_i \|c_i - C(r)\|_2^2\right) . \qquad (10)$$

The background entropy loss regularizes the rendered background probability, $T_{K+1}$ in Eq. (2), to concentrate on either foreground or background:

$$\mathcal{L}_{\text{bg}} = -T_{K+1} \log(T_{K+1}) - (1 - T_{K+1}) \log(1 - T_{K+1}) . \qquad (11)$$

Finally, the overall training objective of the coarse stage is

$$\mathcal{L}^{(\text{c})} = w_{\text{photo}}^{(\text{c})} \cdot \mathcal{L}_{\text{photo}} + w_{\text{pt\_rgb}}^{(\text{c})} \cdot \mathcal{L}_{\text{pt\_rgb}} + w_{\text{bg}}^{(\text{c})} \cdot \mathcal{L}_{\text{bg}} , \qquad (12)$$

where $w_{\text{photo}}^{(\text{c})}, w_{\text{pt\_rgb}}^{(\text{c})}, w_{\text{bg}}^{(\text{c})}$ are the loss weights.

## 5.2. Fine detail reconstruction

Given the optimized coarse geometry $V^{(\text{density})(\text{c})}$ in Sec. 5.1, we now can focus on a smaller subspace to reconstruct the surface details and view-dependent effects. The optimized $V^{(\text{density})(\text{c})}$ is frozen in this stage.

**Fine scene representation.** In the fine stage, we use a higher-resolution density voxel grid $V^{(\text{density})(\text{f})} \in \mathbb{R}^{1 \times N_x^{(\text{f})} \times N_y^{(\text{f})} \times N_z^{(\text{f})}}$ with post-activated interpolation (Eq. (6c)). Note that, alternatively, it is also possible to use a more advanced data structure [18, 30, 67] to refine the voxel grid based on the current $V^{(\text{density})(\text{c})}$ but we leave that for future work. To model view-dependent color emission, we opt to

use an explicit-implicit hybrid representation as we find in our prior experiments that an explicit representation tends to produce worse results, and an implicit representation entails a slower training speed. Our hybrid representation comprises *i)* a feature voxel grid $V^{(\text{feat})(\text{f})} \in \mathbb{R}^{D \times N_x^{(\text{f})} \times N_y^{(\text{f})} \times N_z^{(\text{f})}}$, where $D$ is a hyperparameter for feature-space dimension, and *ii)* a shallow MLP parameteriszed by $\Theta$. Finally, queries of 3D points $\boldsymbol{x}$ and viewing-direction $\boldsymbol{d}$ are performed by

$$\ddot{\sigma}^{(\text{f})} = \text{interp}\left(\boldsymbol{x}, \boldsymbol{V}^{(\text{density})(\text{f})}\right) , \tag{13a}$$

$$\boldsymbol{c}^{(\text{f})} = \text{MLP}_{\Theta}^{(\text{rgb})}\left(\text{interp}(\boldsymbol{x}, \boldsymbol{V}^{(\text{feat})(\text{f})}), \boldsymbol{x}, \boldsymbol{d}\right) , \tag{13b}$$

where $\boldsymbol{c}^{(\text{f})} \in \mathbb{R}^3$ is the view-dependent color emission and $\ddot{\sigma}^{(\text{f})} \in \mathbb{R}$ is the raw volume density in the fine stage. Positional embedding [37] is applied on $\boldsymbol{x}, \boldsymbol{d}$ for the $\text{MLP}_{\Theta}^{(\text{rgb})}$.

**Known free space and unknown space.** A query point is in the known free space if the post-activated alpha value from the optimized $V^{(\text{density})(\text{c})}$ is less than the threshold $\tau^{(\text{c})}$. Otherwise, we say the query point is in the unknown space.

**Fine voxels allocation.** We densely query $V^{(\text{density})(\text{c})}$ to find a BBox tightly enclosing the unknown space, where $L_x^{(\text{f})}, L_y^{(\text{f})}, L_z^{(\text{f})}$ are the lengths of the BBox. The only hyperparameter is the expected total number of voxels $M^{(\text{f})}$. The voxel size $s^{(\text{f})}$ and the grid dimensions $N_x^{(\text{f})}, N_y^{(\text{f})}, N_z^{(\text{f})}$ can then be derived automatically from $M^{(\text{f})}$ as per Sec. 5.1.

**Progressive scaling.** Inspired by NSVF [30], we progressively scale our voxel grid $V^{(\text{density})(\text{f})}$ and $V^{(\text{feat})(\text{f})}$. Let pg_ckpt be the set of checkpoint steps. The initial number of voxels is set to $\lfloor M^{(\text{f})}/2^{|\text{pg\_ckpt}|} \rfloor$. When reaching the training step in pg_ckpt, we double the number of voxels such that the number of voxels after the last checkpoint is $M^{(\text{f})}$; the voxel size $s^{(\text{f})}$ and the grid dimensions $N_x^{(\text{f})}, N_y^{(\text{f})}, N_z^{(\text{f})}$ are updated accordingly. Scaling our scene representation is much simpler. At each checkpoint, we resize our voxel grids, $V^{(\text{density})(\text{f})}$ and $V^{(\text{feat})(\text{f})}$, by trilinear interpolation.

**Fine-stage points sampling.** The points sampling strategy is similar to Eq. (8) with some modifications. We first filter out rays that do not intersect with the known free space. For each ray, we adjust the near- and far-bound, $t^{(\text{near})}$ and $t^{(\text{far})}$, to the two endpoints of the ray-box intersection. We do not adjust $t^{(\text{near})}$ if $\boldsymbol{x}_0$ is already inside the BBox.

**Free space skipping.** Querying $V^{(\text{density})(\text{c})}$ (Eq. (7a)) is faster than querying $V^{(\text{density})(\text{f})}$ (Eq. (13a)); querying for view-dependent colors (Eq. (13b)) is the slowest. We improve fine-stage efficiency by free space skipping in both training and testing. First, we skip sampled points that are in the known free space by checking the optimized $V^{(\text{density})(\text{c})}$ (Eq. (7a)). Second, we further skip sampled points in unknown space with low activated alpha value (threshold at $\tau^{(\text{f})}$) by querying $V^{(\text{density})(\text{f})}$ (Eq. (13a)).

**Optimizing the fine scene representation.** The training objective is similar to Eq. (12), which is

$$\mathcal{L}^{(\text{f})} = w_{\text{photo}}^{(\text{f})} \cdot \mathcal{L}_{\text{photo}} + w_{\text{pt\_rgb}}^{(\text{f})} \cdot \mathcal{L}_{\text{pt\_rgb}} + w_{\text{bg}}^{(\text{f})} \cdot \mathcal{L}_{\text{bg}} , \tag{14}$$

where $w_{\text{photo}}^{(\text{f})}, w_{\text{pt\_rgb}}^{(\text{f})}, w_{\text{bg}}^{(\text{f})}$ are the loss weights in fine stage.

# 6. Experiments

## 6.1. Datasets

We evaluate our approach on five inward-facing datasets. **Synthetic-NeRF** [37] contains eight objects with realistic images synthesized by NeRF. **Synthetic-NSVF** [30] contains another eight objects synthesized by NSVF. Strictly following NeRF's and NSVF's setups, we set the image resolution to $800 \times 800$ pixels and let each scene have 100 views for training and 200 views for testing. **BlendedMVS** [66] is a synthetic MVS dataset that has realistic ambient lighting from real image blending. We use a subset of four objects provided by NSVF. The image resolution is $768 \times 576$ pixels, and one-eighth of the images are for testing. **Tanks&Temples** [21] is a real-world dataset. We use a subset of five scenes provided by NSVF, each containing views captured by an inward-facing camera circling the scene. The image resolution is $1920 \times 1080$ pixels, and one-eighth of the images are for testing. **DeepVoxels** [49] dataset contains four simple Lambertian objects. The image resolutions are $512 \times 512$, and each scene has 479 views for training and 1000 views for testing.

## 6.2. Implementation details

We choose the same hyperparameters generally for all scenes. The expected numbers of voxels are set to $M^{(\text{c})} = 100^3$ and $M^{(\text{f})} = 160^3$ in coarse and fine stages. The activated alpha values are initialized to be $\alpha^{(\text{init})(\text{c})} = 10^{-6}$ in the coarse stage. We use a higher $\alpha^{(\text{init})(\text{f})} = 10^{-2}$ as the query points are concentrated on the optimized coarse geometry in the fine stage. The points sampling step sizes are set to half of the voxel sizes, *i.e.*, $\delta^{(\text{c})} = 0.5 \cdot s^{(\text{c})}$ and $\delta^{(\text{f})} = 0.5 \cdot s^{(\text{f})}$. The shallow MLP layer comprises two hidden layers with 128 channels. We use the Adam optimizer [20] with a batch size of 8,192 rays to optimize the coarse and fine scene representations for 10k and 20k iterations. The base learning rates are 0.1 for all voxel grids and $10^{-3}$ for the shallow MLP. The exponential learning rate decay is applied. See supplementary material for detailed hyperparameter setups.

## 6.3. Comparisons

**Quantitative evaluation on the synthesized novel view.** We first quantitatively compare the novel view synthesis results in Tab. 1. PSNR, SSIM [61], and LPIPS [70] are employed as evaluation metrics. We outperform the original NeRF [37] and the improved JaxNeRF [9] re-implementation. Besides, our results are also comparable

| Methods | Synthetic-NeRF | | | Synthetic-NSVF | | | BlendedMVS | | | Tanks and Temples | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| SRN [50] | 22.26 | 0.846 | 0.170$^{\text{vgg}}$ | 24.33 | 0.882 | 0.141$^{\text{alex}}$ | 20.51 | 0.770 | 0.294$^{\text{alex}}$ | 24.10 | 0.847 | 0.251$^{\text{alex}}$ |
| NV [32] | 26.05 | 0.893 | 0.160$^{\text{vgg}}$ | 25.83 | 0.892 | 0.124$^{\text{alex}}$ | 23.03 | 0.793 | 0.243$^{\text{alex}}$ | 23.70 | 0.834 | 0.260$^{\text{alex}}$ |
| NeRF [37] | 31.01 | 0.947 | 0.081$^{\text{vgg}}$ | 30.81 | 0.952 | 0.043$^{\text{alex}}$ | 24.15 | 0.828 | 0.192$^{\text{alex}}$ | 25.78 | 0.864 | 0.198$^{\text{alex}}$ |
| *Improved visual quality from NeRF* | | | | | | | | | | | | |
| JaxNeRF [9] | 31.69 | 0.953 | 0.068$^{\text{vgg}}$ | - | - | - | - | - | - | 27.94 | 0.904 | 0.168$^{\text{vgg}}$ |
| JaxNeRF+ [9] | 33.00 | 0.962 | 0.038 | - | - | - | - | - | - | - | - | - |
| Mip-NeRF [1] | 33.09 | 0.961 | 0.043$^{\text{vgg}}$ | - | - | - | - | - | - | - | - | - |
| *Improved test-time rendering speed (and visual quality) from NeRF* | | | | | | | | | | | | |
| AutoInt [29] | 25.55 | 0.911 | 0.170 | - | - | - | - | - | - | - | - | - |
| FastNeRF [15] | 29.97 | 0.941 | 0.053 | - | - | - | - | - | - | - | - | - |
| SNeRG [18] | 30.38 | 0.950 | 0.050 | - | - | - | - | - | - | - | - | - |
| KiloNeRF [44] | 31.00 | 0.95 | 0.03 | 33.37 | 0.97 | 0.02 | 27.39 | 0.92 | 0.06 | 28.41 | 0.91 | 0.09 |
| PlenOctrees [67] | 31.71 | 0.958 | 0.053$^{\text{vgg}}$ | - | - | - | - | - | - | 27.99 | 0.917 | 0.131$^{\text{vgg}}$ |
| NSVF [30] | 31.75 | 0.953 | 0.047$^{\text{alex}}$ | 35.18 | 0.979 | 0.015$^{\text{alex}}$ | 26.89 | 0.898 | 0.114$^{\text{alex}}$ | 28.48 | 0.901 | 0.155$^{\text{alex}}$ |
| *Improved convergence speed, test-time rendering speed, and visual quality from NeRF* | | | | | | | | | | | | |
| ours | 31.95 | 0.957 | 0.053$^{\text{vgg}}$ 0.035$^{\text{alex}}$ | 35.08 | 0.975 | 0.033$^{\text{vgg}}$ 0.019$^{\text{alex}}$ | 28.02 | 0.922 | 0.101$^{\text{vgg}}$ 0.075$^{\text{alex}}$ | 28.41 | 0.911 | 0.155$^{\text{vgg}}$ 0.148$^{\text{alex}}$ |

* The superscript denotes the pre-trained models used in LPIPS. The gray numbers indicate that the code is unavailable or has a unconventional LPIPS implementation.

Table 1. **Quantitative comparisons for novel view synthesis.** Our method excels in convergence speed, *i.e.*, 15 minutes per scene compared to many hours or days per scene using other methods. Besides, our rendering quality is better than the original NeRF [37] and the improved JaxNeRF [9] on the four datasets under all metrics. We also show comparable results to most of the recent methods.

to most of the recent methods, except JaxNeRF+ [9] and Mip-NeRF [1]. Moreover, our per-scene optimization only takes about 15 minutes, while all the methods after NeRF in Tab. 1 need quite a few hours per scene. We defer detail comparisons on the much simpler DeepVoxels [49] dataset to supplementary material, where we achieve 45.83 averaged PSNR and outperform NeRF's 40.15 and IBRNet's 42.93.

**Training time comparisons.** The key merit of our work is the significant improvement in convergence speed with NeRF-comparable quality. In Tab. 2, we show a training time comparison. We also show GPU specifications after each reported time as it is the main factor affecting run-time.

NeRF [37] with a more powerful GPU needs 1–2 days per scene to achieve 31.01 PSNR, while our method achieves a superior 31.95 PSNR in about 15 minutes per scene. MVS-NeRF [6], IBRNet [60], and NeuRay [31] also show less per-scene training time than NeRF but with the additional cost to run a generalizable cross-scene pre-training. MVS-NeRF [6], after pre-training, optimizes a scene in 15 minutes as well, but the PSNR is degraded to 28.14. IBRNet [60] shows worse PSNR and longer training time than ours. Neu-Ray [31] originally reports time in lower-resolution (NeuRay-Lo) setup, and we receive the training time of the high-resolution (NeuRay-Hi) setup from the authors. NeuRay-Hi achieves a better 32.42 PSNR but requires 23 hours to train. For the early-stopped NeuRay-Hi, unfortunately, only its training time is retained (early-stopped NeuRay-Lo achieves NeRF-similar PSNR). Even with early stopping, NeuRay-Hi still needs 70 minutes to train, while we only need 15 min-

| Methods | PSNR↑ | generalizable pre-training | per-scene optimization |
|---|---|---|---|
| NeRF [37] | 31.01 | no need | 1–2 days (V100) |
| MVSNeRF [6] | 27.21 | 30 hrs (2080Ti) | 15 mins (2080Ti) |
| IBRNet [60] | 28.14 | 1 day (8×V100) | 6 hrs (V100) |
| NeuRay [31]† | 32.42 | 2 days (2080Ti) | 23 hrs (2080Ti) |
| Mip-NeRF [1]‡ | 30.85 | no need | 6 hrs (2080Ti) |
| ours | 31.95 | no need | 15 mins (2080Ti) |

† Use external depth information.

‡ Our reproduction with early stopping on our machine.

Table 2. **Training time comparisons.** We take the training time and GPU specifications reported in previous works directly. A V100 GPU can run faster and has more storage than a 2080Ti GPU. We include mean PSNR on the eight scenes in the *Synthetic-NeRF* [37] dataset as an assessment of rendering quality. Our method achieves good PSNR in a significantly less per-scene optimization time.

utes and do not rely on generalizable pre-training or external depth information. Mip-NeRF [1] has similar run-time to NeRF but with much better PSNRs, which also signifies using less training time to achieve NeRF's PSNR. We train early-stopped Mip-NeRFs on our machine and show the averaged PSNR and training time. The early-stopped Mip-NeRF achieves 30.85 PSNR after 6 hours of training, while we can achieve 31.95 PSNR in just 15 minutes.

**Rendering speed comparisons.** Improving test-time rendering speed is not the main focus of this work, but we still achieve $\sim 45\times$ speedups from NeRF—0.64 seconds versus 29 seconds per $800 \times 800$ image on our machine.
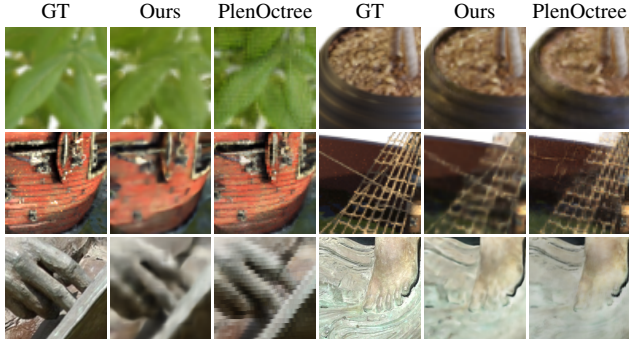
GT    Ours    PlenOctree    GT    Ours    PlenOctree

Figure 5. **Qualitative comparisons on the challenging parts.** *Top*: On *ficus* scene, we do not show blocking artifacts as PlenOctree and recover the pot better. *Middle*: We produce blurrier results on *ship*'s body and rigging, but we do not have the background artifacts. *Bottom*: On real-world captured *Ignatius*, we show better quality without blocking artifacts (left) and recover the color tone better (right). See supplementary material for more visualizations.

**Qualitative comparison.** Fig. 5 shows our rendering results on the challenging parts and compare them with the results (better than NeRF's) provided by PlenOctrees [67].

## 6.4. Ablation studies

We mainly validate the effectiveness of the two proposed techniques—post-activation and the imposed priors—that enable voxel grids to model scene geometry with NeRF-comparable quality. We subsample two scenes for each dataset. See supplementary material for more detail and additional ablation studies on the number of voxels, point-sampling step size, progressive scaling, free space skipping, view-dependent colors modeling, and the losses.

**Effectiveness of the post-activation.** We show in Sec. 4 that the proposed post-activated trilinear interpolation enables the discretized grid to model sharper surfaces. In Tab. 3, we compare the effectiveness of post-activation in scene reconstruction for novel view synthesis. Our grid in the fine stage consists of only $160^3$ voxels, where nearest-neighbor interpolation results in worse quality than trilinear interpolation. The proposed post-activation can improve the results further compared to pre- and in-activation. We find that we gain less in the real-world captured *BlendedMVS* and *Tanks and Temples* datasets. The intuitive reason is that real-world data introduces more uncertainty (*e.g.*, inconsistent lightning, SfM error), which results in multi-view inconsistent and blurrier surfaces. Thus, the advantage is lessened for scene representations that can model sharper surfaces. We speculate that resolving the uncertainty in future work can increase the gain of the proposed post-activation.

**Effectiveness of the imposed priors.** As discussed in Sec. 5.1, it is crucial to initialize the voxel grid with low density to avoid suboptimal geometry. The hyperparameter

| Interp. | | Syn.-NeRF | | Syn.-NSVF | | BlendedMVS | | T&T | |
|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ |
| Nearest | | 28.61 | -2.77 | 28.86 | -6.22 | 25.49 | -2.48 | 26.39 | -1.27 |
| Tri. pre- | | 30.84 | -0.55 | 32.66 | -2.41 | 27.39 | -0.58 | 27.44 | -0.21 |
| Tri. in- | | 29.91 | -1.48 | 32.42 | -2.66 | 27.29 | -0.68 | 27.52 | -0.13 |
| Tri. post- | | **31.39** | - | **35.08** | - | **27.97** | - | **27.66** | - |

Table 3. **Effectiveness of the post-activation.** Geometry modeling with density voxel grid can achieve better PSNRs by using the proposed post-activated trilinear interpolation.

| $\alpha^{(init)(c)}$ | View. lr. | Syn.-NeRF | | Syn.-NSVF | | BlendedMVS | | T&T | |
|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ |
| - | ✓ | 28.88 | -2.51 | 25.12 | -9.96 | 22.17 | -5.79 | 25.33 | -2.33 |
| $10^{-3}$ | ✓ | 30.96 | -0.42 | 27.24 | -7.84 | 23.17 | -4.79 | 26.04 | -1.61 |
| $10^{-4}$ | ✓ | 31.29 | -0.09 | 31.05 | -4.03 | 26.09 | -1.88 | 27.60 | -0.05 |
| $10^{-5}$ | ✓ | 31.41 | +0.02 | 35.04 | -0.04 | 27.36 | -0.61 | 27.63 | -0.02 |
| $10^{-6}$ | | 31.40 | +0.01 | 35.03 | -0.04 | 27.37 | -0.60 | 27.59 | -0.07 |
| $10^{-7}$ | ✓ | 31.36 | -0.02 | 35.03 | -0.05 | 27.73 | -0.23 | 27.59 | -0.06 |
| $10^{-6}$ | ✓ | **31.39** | - | **35.08** | - | **27.97** | - | **27.66** | - |



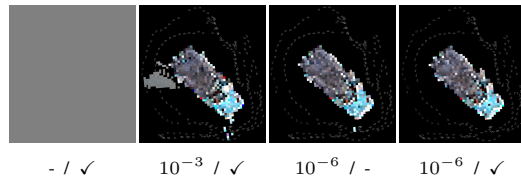- / ✓        $10^{-3}$ / ✓        $10^{-6}$ / -        $10^{-6}$ / ✓

Table 4. **Effectiveness of the imposed priors.** We compare our different settings in the coarse geometry search. *Top*: We show their impacts on the final PSNRs after the fine stage reconstruction. *Bottom*: We visualize the allocated voxels by coarse geometry search on the *Truck* scene. Overall, low-density initialization is essential; using $\alpha^{(init)(c)} = 10^{-6}$ and view-count-based learning rate generally achieves cleaner voxels allocation in the coarse stage and better PSNR after the fine stage.

$\alpha^{(init)(c)}$ controls the initial activated alpha values via Eq. (9). In Tab. 4, we compare the quality with different $\alpha^{(init)(c)}$ and the view-count-based learning rate. Without the low-density initialization, the quality drops severely for all the scenes. When $\alpha^{(init)(c)} = 10^{-7}$, we have to train the coarse stage of some scenes for more iterations. The effective range of $\alpha^{(init)(c)}$ is scene-dependent. We find $\alpha^{(init)(c)} = 10^{-6}$ generally works well on all the scenes in this work. Finally, using a view-count-based learning rate can further improve the results and allocate noiseless voxels in the coarse stage.

## 7. Conclusion

Our method directly optimizes the voxel grid and achieves super-fast convergence in per-scene optimization with NeRF-comparable quality—reducing training time from many hours to 15 minutes. However, we do not deal with the unbounded or forward-facing scenes, while we believe our method can be a stepping stone toward fast convergence in such scenarios. We hope our method can boost the progress of NeRF-based scene reconstruction and its applications.

# Supplementary material

In Sec. H, we give more implementation details. In Sec. I, we present additional ablation studies for the important hyperparameters that affect the quality and convergence speed. In Sec. J, we show detailed results of our main ablation studies presented in the main paper. In Sec. K, we report our training time breakdown of the coarse and fine stage on each scene. In Sec. L, we show the detailed quantitative and qualitative results on each scene. Finally, we derive for the low-density initialization in Sec. M, and we prove that the proposed post-activated trilinear interpolation can produce sharp linear surfaces in Sec. N.

## H. Additional implementation details

We choose the same hyperparameters generally for all scenes in the five datasets.

In the coarse stage, the expected number of voxels is $M^{(c)} = 100^3$. The activated alpha values are initialized to be $\alpha^{(init)(c)} = 10^{-6}$.

In the fine stage, we use $M^{(f)} = 160^3$ voxels. We use a higher $\alpha^{(init)(f)} = 0.01$ as the query points in the known free space are skipped. The shallow MLP layer comprises two hidden layers with 128 channels. The number of channels of the feature voxel grid $V^{(feat)(f)}$ is set to $D = 12$. The number of frequencies of positional embedding for the query position $x$ and the viewing direction $d$ are $k_x = 5, k_d = 4$. We progressively scale the fine-stage voxel grid at the checkpoint pg_ckpt being chosen as the 1000th, 2000th, and 3000th training steps. We use threshold $\tau^{(c)} = 10^{-3}$ to define the known free space and $\tau^{(f)} = 10^{-4}$ to skip the low-density query points in unknown space.

We set the sampling step sizes to be half of the voxel sizes, i.e., $\delta^{(c)} = 0.5s^{(c)}$ and $\delta^{(f)} = 0.5s^{(f)}$. During training, we randomly shift the query points on a ray by $\left(p \cdot \delta^{(\cdot)} \cdot d/\|d\|^2\right)$, where $p$ is sampled uniformly in $[0, 1]$, and $d$ is the viewing direction of the ray.

The loss weights are set as follows: $w_{photo}^{(c)} = 1$, $w_{pt\_rgb}^{(c)} = 10^{-1}$, $w_{bg}^{(c)} = 10^{-2}$, $w_{photo}^{(f)} = 1$, $w_{pt\_rgb}^{(f)} = 10^{-2}$, and $w_{bg}^{(f)} = 10^{-3}$. We use the Adam optimizer with a batch size of 8,192 rays to optimize the coarse and fine scene representations for 10k and 20k iterations. The base learning rates are 0.1 for all voxel grids and $10^{-3}$ for the MLP. The exponential learning rate decay is applied such that the learning rates are downscaled by 0.1 at 20k iterations.

## I. Additional ablation experiments

We use the *Robot* scene to conduct the additional ablation experiments. In all the additional ablation experiments, the fine stage are early stopped at 10k iterations. The reported training times are measured on our machine with a single RTX2080 Ti GPU.

**Effect of the number of voxels.** The hyperparameters $M^{(c)}$ and $M^{(f)}$ define the numbers of voxels in the coarse stage and the fine stage. In Tab. I1, we show the effect of different $M^{(\cdot)}$ setups on the final PSNRs and the training times. We can use more voxels for better quality or use fewer voxels for faster convergence speed.

| $M^{(c)}$ | $M^{(f)}$ | PSNR↑ | Training minutes↓ |
|---|---|---|---|
| $100^3$ | $100^3$ | 32.65 | **5.1** |
| $100^3$ | $136^3$ | 34.57 | 7.0 |
| $100^3$ | $160^3$ | **35.54** | 7.3 |

(a) Using fewer voxels in the fine stage results in inferior qualities but faster training speed.

| $M^{(c)}$ | $M^{(f)}$ | PSNR↑ | Training minutes↓ |
|---|---|---|---|
| $64^3$ | $160^3$ | 34.91 | **6.1** |
| $100^3$ | $160^3$ | 35.54 | 7.3 |
| $136^3$ | $160^3$ | **35.84** | 9.3 |

(b) The number of voxels in the coarse stage can also affect final results as the fine stage relies on the coarse geometry optimized in the coarse stage.

Table I1. Effect of the number of voxels in the coarse stage ($M^{(c)}$) and the fine stage ($M^{(f)}$). Using more voxels can improve quality with the cost of more computation and training time.

**Effect of the point sampling step size.** In our design, the voxel sizes $s^{(\cdot)}$ are automatically derived from the hyperparameters of the number of voxels $M^{(\cdot)}$. We set the step-size-to-voxel-size ratio instead of directly assigning the step size for points sampling on rays, i.e., $\delta^{(\cdot)} = 0.5s^{(\cdot)}$ means that the step size is half a voxel size. We show in Tab. I2 that the step sizes are also important hyperparameters for the speed-quality tradeoff. We can also improve the rendering quality of a trained model by using a finer step size. For the sake of simplicity, we leave future work to explore the effect of different hierarchical point sampling strategies [1,37,39].

**Effect of the progressive scaling in the fine stage.** In the fine-stage training, we scale the resolution of our voxel grids progressively. We show in Tab. I3 that such a strategy can improve the training efficiency with a slightly better rendering quality.

**Effect of the free space skipping in the fine stage.** There are three models in the fine stage—*i)* the optimized and frozen coarse density voxel grid $V^{(density)(c)}$, *ii)* the finer density voxel grid $V^{(density)(f)}$, and *iii)* the model for the view-dependent color emission. Querying the optimized $V^{(density)(c)}$ is more efficient than querying the finer density $V^{(density)(f)}$; querying the view-dependent color emission is the slowest and computationally expensive. For each query point, we first query $V^{(density)(c)}$ and ignore the query point if the activated alpha is below a threshold $\tau^{(c)}$ (i.e., the point is in the known free space). As the training progresses, we can filter out more points that are in the free space defined by the

| Step size | | PSNR↑ | sec/frame↓ |
|---|---|---|---|
| Training | Testing | | |
| $\delta^{(\cdot)} = 0.50s^{(\cdot)}$ | $\delta^{(\cdot)} = 0.50s^{(\cdot)}$ | 35.54 | **0.64** |
| $\delta^{(\cdot)} = 0.50s^{(\cdot)}$ | $\delta^{(\cdot)} = 0.25s^{(\cdot)}$ | **35.72** | 1.15 |

(a) We can improve quality by using a finer step size in test-time with the cost of slower rendering speed. The two results share a trained model and are only different in the testing step size.

| Step size | PSNR↑ | Training minutes↓ |
|---|---|---|
| $\delta^{(\cdot)} = 1.00s^{(\cdot)}$ | 34.17 | **5.2** |
| $\delta^{(\cdot)} = 0.75s^{(\cdot)}$ | 34.96 | 5.9 |
| $\delta^{(\cdot)} = 0.50s^{(\cdot)}$ | 35.54 | 7.3 |
| $\delta^{(\cdot)} = 0.25s^{(\cdot)}$ | **36.01** | 11.5 |

(b) Training with finer step size can improve results with the cost of more computation and training time.

Table I2. Effect of the step size in points sampling. The step size $\delta^{(\cdot)}$ is relative to the voxel size $s^{(\cdot)}$, and we apply the same step-size-to-voxel-size ratio in our coarse and fine stage. We can achieve better quality by using a smaller (finer) step size with the cost of more computation.

| Progressive scaling | PSNR↑ | Training minutes↓ |
|---|---|---|
|  | 35.50 | 11.7 |
| ✓ | **35.54** | **7.3** |

Table I3. Scaling the grid resolutions progressively in the fine stage can improve training efficiency with a slightly better quality.

"sculpted" $V^{(\text{density})(f)}$. Specifically, we ignore a query point if the activated alpha from $V^{(\text{density})(f)}$ is below a threshold $\tau^{(f)}$. Thus, we only query for the view-dependent color if the query point passes the two criteria.

In Tab. I4, we show the effect of the free space skipping strategy. In the case of $\tau^{(c)} = 0$, the coarse voxel grid $V^{(\text{density})(c)}$ is only used to determine a tighten BBox enclosing the scene. We finally run out of memory (OOM) at the fine stage if no skipping rule is applied ($\tau^{(c)} = 0, \tau^{(f)} = 0$). When only $\tau^{(f)}$ is applied, we can still sculpt out many irrelevant query points during the progressive scaling and save the memory before our voxel grids are scaled to the finest resolution. The rendering quality is degraded when $\tau^{(c)} = 0$, which suggests that the imposed priors in the coarse stage are helpful to the final quality (we cannot apply the low-density initialization and the free space skipping at the same time). Finally, we achieve the best rendering quality and convergence speed when both $\tau^{(c)}$ and $\tau^{(f)}$ are applied.

**Ablation studies for the view-dependent color modeling.** Our hybrid representation for the view-dependent color emissions in the fine stage comprises a feature voxel grid and a shallow MLP. In Tab. I5, we show different strategies and hyperparameters to model view-dependent colors. Applying only the shallow MLP results in worse outputs as our MLP

| Free space skipping | | PSNR↑ | Training minutes↓ |
|---|---|---|---|
| $\tau^{(c)}$ | $\tau^{(f)}$ | | |
| 0 | 0 | - | OOM |
| 0 | $10^{-4}$ | 34.89 | 7.8 |
| $10^{-3}$ | 0 | **35.54** | 15.6 |
| $10^{-3}$ | $10^{-4}$ | **35.54** | **7.3** |

Table I4. Effect of the free space skipping strategy in the fine stage training.

is much "shallower" than the MLP in NeRF—our two layers with 128 channels versus NeRF's eight layers with 256 channels. Using only the voxel grid to model view-invariant diffused colors can converge in much less time, but the rendering quality is degraded. In NeRF-based scene reconstruction, spherical harmonic representation [31, 67], learnable basis [63], or deferred rendering [18] are also shown to be effective, but we leave the explorations for future work. We also compare the number of dimensions $D$ of our feature grid, the number of layers, and the number of channels of the shallow MLP. Using a larger model leads to better rendering quality with the cost of more computation.

| Colors rep. in fine stage | | | PSNR↑ | Training minutes↓ |
|---|---|---|---|---|
| implicit (shallow MLP) | | | 26.57 | 6.2 |
| explicit (diffused only) | | | 32.15 | **2.9** |
| hybrid | | | | |
| $D$ | MLP layers | MLP ch. | | |
| 12 | 1 | 128 | 35.11 | 7.0 |
| 12 | 2 | 128 | 35.54 | 7.3 |
| 12 | 3 | 128 | 35.70 | 7.8 |
| 12 | 2 | 64 | 35.35 | 7.0 |
| 12 | 2 | 192 | 35.71 | 7.7 |
| 24 | 2 | 128 | **35.90** | 8.8 |

Table I5. Different representations and hyperparameters for modeling the view-dependent color emissions.

**Effect of the auxiliary losses.** We show in Tab. I6 that incorporating the per-point rgb loss and the background entropy loss added to the main photometric loss can improve our results. The per-point rgb loss supervises the color emission of each query points directly instead of the accumulated color. The background entropy loss enforces the rendered background probability to concentrate on the background or foreground. We achieve the best results when adding the two auxiliary losses.

## J. Main ablation studies details

In Tab. J1, we show detailed results for the ablation experiments presented in the main paper. We subsample two scenes for each dataset to conduct the main ablation experiments—

| Per-point rgb loss | Background entropy loss | PSNR↑ |
|:---:|:---:|:---:|
| | | 34.84 |
| ✓ | | 35.37 |
| | ✓ | 35.53 |
| ✓ | ✓ | **35.54** |

Table I6. Effect of the auxiliary losses.

*Materials* and *Mic* from Synthetic-NeRF dataset; *Robot* and *Lifestyle* from Synthetic-NSVF dataset; *Character* and *Statues* from BlendedMVS dataset, and *Ignatius* and *Truck* from Tanks and Temples dataset.

We show that the proposed post-activation significantly improves our quality. The low-density initialization is shown to be essential to our method and is an important hyperparameter. The view-count-based learning rate can slightly improve the results, and the PSNR without the view-count prior can degrade up to $-1.22$ in the worst case.

## K. Additional training time details

### K.1. Training time comparisons

Mip-NeRF [1] requires similar run-time as NeRF [37] but achieves much better quality. As a side benefit, Mip-NeRF can attain NeRF's PSNR in less time. To compare the convergence speed with Mip-NeRF, we use the code open-sourced by the authors (https://github.com/google/mipnerf). We re-train early-stopped Mip-NeRFs on our machine with only a single RTX2080 Ti GPU. The comparison is shown in Tab. K1. The early-stopped Mip-NeRF that is trained for 6 hours achieves an average PSNR of 30.85. On the other hand, our method is only optimized for about 15 minutes and achieves an average PSNR of 31.95.

### K.2. The detailed training time of our model

We report our training times in detail on each scene and each stage in Tab. K2. The coarse stage optimization accounts for about $15-20\%$ of the overall training time. The per-scene training times are about less than 15 minutes, except for the *Tanks&Temples* [21] dataset, which takes just a few more minutes.

## L. Per-scene analysis

### L.1. Per-scene quantitative results

We report the per-scene quantitative comparisons in Tab. L1 for **Synthetic-NeRF** [37] dataset, Tab. L2 for **Synthetic-NSVF** [30] dataset, Tab. L3 for **Blended-MVS** [66] dataset, Tab. L4 for **Tanks&Temples** [21] dataset, and Tab. L5 for **DeepVoxels** [49] dataset. Our method achieves comparable results to most of the recent methods, except the Mip-NeRF [1] and JaxNeRF+ [9]. Moreover, all

the methods after NeRF on the tables take quite a few hours to train for each scene, while our method only takes about 15 minutes per-scene optimization time as reported in Tab. K2.

### L.2. Per-scene qualitative results

We provide more qualitative results in Fig. L2 for **Synthetic-NSVF** [30] dataset, Fig. L3 for **Blended-MVS** [66] dataset, and Fig. L5 for **DeepVoxels** [49] dataset. In Fig. L1, we compare our results with the results provided by JaxNeRF [9] and PlenOctrees [67] on **Synthetic-NeRF** [37] dataset. Both JaxNeRF and PlenOctrees are stronger baselines than the original NeRF [37]. In Fig. L4, we also compare our results with PlenOctrees on the **Tanks&Temples** [21] dataset. In the qualitative comparisons, there is no consistent superior method across all the scenes, which matches the results in quantitative comparisons.

| Interp. | $\alpha^{(init)(c)}$ | View. lr. | Overall | | | | Synthetic-NeRF | | | | Synthetic-NSVF | | | | BlendedMVS | | | | Tanks and Temples | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | *Materials* | | *Mic* | | *Robot* | | *Lifestyle* | | *Character* | | *Statues* | | *Ignatius* | | *Truck* | |
| | | | PSNR↑ | $\Delta_{avg}$ | $\Delta_{worst}$ | $\Delta_{best}$ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ | PSNR↑ | Δ |
| post-act. | $10^{-6}$ | ✓ | 30.52 | - | - | - | 29.57 | - | 33.20 | - | 36.36 | - | 33.79 | - | 30.31 | - | 25.62 | - | 28.16 | - | 27.15 | |
| nearest | $10^{-6}$ | ✓ | 27.33 | -3.19 | -7.50 | -0.34 | 28.17 | -1.40 | 29.05 | -4.15 | 28.86 | -7.50 | 28.85 | -4.94 | 27.28 | -3.03 | 23.69 | -1.93 | 27.82 | -0.34 | 24.95 | -2.20 |
| pre-act. | $10^{-6}$ | ✓ | 29.58 | -0.94 | -2.88 | -0.14 | 29.24 | -0.33 | 32.43 | -0.77 | 33.48 | -2.88 | 31.85 | -1.94 | 29.73 | -0.58 | 25.04 | -0.58 | 28.02 | -0.14 | 26.86 | -0.29 |
| in-act. | $10^{-6}$ | ✓ | 29.28 | -1.24 | -3.30 | 0.00 | 27.34 | -2.23 | 32.47 | -0.73 | 33.06 | -3.30 | 31.77 | -2.02 | 29.74 | -0.57 | 24.84 | -0.78 | 28.16 | 0.00 | 26.89 | -0.26 |
| post-act. | - | ✓ | 25.37 | -5.15 | -9.98 | -1.61 | 27.40 | -2.17 | 30.35 | -2.85 | 26.43 | -9.93 | 23.81 | -9.98 | 24.70 | -5.61 | 19.65 | -5.97 | 26.55 | -1.61 | 24.10 | -3.05 |
| | $10^{-3}$ | ✓ | 26.85 | -3.67 | -8.40 | -0.23 | 28.96 | -0.61 | 32.97 | -0.23 | 29.09 | -7.27 | 25.39 | -8.40 | 25.12 | -5.19 | 21.22 | -4.40 | 27.23 | -0.93 | 24.86 | -2.29 |
| | $10^{-4}$ | ✓ | 29.01 | -1.51 | -4.30 | +0.05 | 29.35 | -0.22 | 33.23 | +0.03 | 32.06 | -4.30 | 30.04 | -3.75 | 28.47 | -1.84 | 23.70 | -1.92 | 28.21 | +0.05 | 26.99 | -0.16 |
| | $10^{-5}$ | ✓ | 30.36 | -0.16 | -1.22 | +0.03 | 29.60 | +0.03 | 33.22 | +0.02 | 36.32 | -0.04 | 33.75 | -0.04 | 30.32 | +0.01 | 24.40 | -1.22 | 28.17 | +0.01 | 27.10 | -0.05 |
| | $10^{-6}$ | ✓ | 30.35 | -0.17 | -1.22 | +0.02 | 29.57 | 0.00 | 33.22 | +0.02 | 36.33 | -0.03 | 33.74 | -0.05 | 30.33 | +0.02 | 24.40 | -1.22 | 28.06 | -0.10 | 27.12 | -0.03 |
| | $10^{-7}$ | ✓ | 30.43 | -0.09 | -0.56 | +0.20 | 29.57 | 0.00 | 33.15 | -0.05 | 36.56 | +0.20 | 33.50 | -0.29 | 30.40 | +0.09 | 25.06 | -0.56 | 28.12 | -0.04 | 27.07 | -0.08 |

Table J1. We detail the per-scene results of different ablation setups presented in the main paper. We also show their difference ($\Delta$) to our final setting (the first row). The $\Delta_{worst}$ highlights the worst-case degradation overall scenes, where most of the settings lead to more than 1 dB degradation in the worst case. The $\Delta_{best}$ shows the best case difference, where some settings can slightly improve the results on some scenes. We highlight the top 3 results of each column in gold, silver, and bronze.

| Methods | | Avg. | *Chair* | *Drums* | *Ficus* | *Hotdog* | *Lego* | *Materials* | *Mic* | *Ship* |
|---|---|---|---|---|---|---|---|---|---|---|
| Mip-NeRF [1] (early-stopped) | PSNR↑ | 30.85 | 32.25 | 24.95 | 30.66 | 35.56 | 32.92 | 29.28 | 32.61 | 28.56 |
| | training minutes↓ | 361.2 | 361.1 | 363.0 | 360.1 | 359.0 | 363.7 | 360.7 | 359.3 | 363.0 |
| Ours | PSNR↑ | **31.95** | **34.09** | **25.44** | **32.78** | **36.74** | **34.64** | **29.57** | **33.20** | **29.13** |
| | training minutes↓ | **14.2** | **12.5** | **12.0** | **13.8** | **15.5** | **13.2** | **15.4** | **11.0** | **20.2** |

Table K1. We compare the training time with the early-stopped Mip-NeRF, which is a stronger baseline than NeRF. The training time is measured on our machine with only a single RTX2080 Ti GPU. Our method with about 15 minutes of training time outperforms the early-stopped Mip-NeRF with 6 hours of training by a large margin.

| Stage | Avg. | *Chair* | *Drums* | *Ficus* | *Hotdog* | *Lego* | *Materials* | *Mic* | *Ship* |
|---|---|---|---|---|---|---|---|---|---|
| Coarse | 2.3 | 2.5 | 2.4 | 2.0 | 2.3 | 2.3 | 2.1 | 2.3 | 2.3 |
| Fine | 11.9 | 10.1 | 9.7 | 11.7 | 13.2 | 10.9 | 13.3 | 8.7 | 17.9 |
| All | 14.2 | 12.5 | 12.0 | 13.8 | 15.5 | 13.2 | 15.4 | 11.0 | 20.2 |

(a) Training times (minutes) on the eight scenes of **Synthetic-NeRF** [37] dataset.

| Stage | Avg. | *Wineholder* | *Steamtrain* | *Toad* | *Robot* | *Bike* | *Palace* | *Spaceship* | *Lifestyle* |
|---|---|---|---|---|---|---|---|---|---|
| Coarse | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 |
| Fine | 10.7 | 9.6 | 12.8 | 9.8 | 9.5 | 10.1 | 11.2 | 12.1 | 10.1 |
| All | 13.2 | 12.1 | 15.3 | 12.3 | 12.0 | 12.6 | 13.7 | 14.6 | 12.6 |

(b) Training times (minutes) on the eight scenes of **Synthetic-NSVF** [30] dataset.

| Stage | Avg. | *Jade* | *Fountain* | *Character* | *Statues* |
|---|---|---|---|---|---|
| Coarse | 2.6 | 2.0 | 2.7 | 3.2 | 2.3 |
| Fine | 11.2 | 11.8 | 10.6 | 11.5 | 11.0 |
| All | 13.8 | 13.8 | 13.3 | 14.7 | 13.2 |

| Stage | Avg. | *Ignatius* | *Truck* | *Barn* | *Cate.* | *Family* |
|---|---|---|---|---|---|---|
| Coarse | 3.6 | 3.4 | 3.6 | 3.8 | 3.6 | 3.4 |
| Fine | 14.2 | 12.4 | 14.3 | 18.2 | 14.1 | 11.9 |
| All | 17.7 | 15.7 | 17.8 | 22.0 | 17.7 | 15.3 |

(c) Training times (minutes) on the four scenes of **BlendedMVS** [66] dataset.  (d) Training times (minutes) on the five scenes of **Tanks&Temples** [21] dataset.

| Stage | Avg. | *Chair* | *Pedestal* | *Cube* | *Vase* |
|---|---|---|---|---|---|
| Coarse | 2.4 | 2.5 | 2.5 | 2.1 | 2.5 |
| Fine | 11.9 | 12.8 | 11.6 | 11.1 | 12.2 |
| All | 14.3 | 15.2 | 14.1 | 13.2 | 14.7 |

(e) Training times (minutes) on the four scenes of **DeepVoxels** [49] dataset.

Table K2. We report the detail of our per-scene optimization time in minutes measured on our machine with a single RTX2080 Ti GPU.

**PSNR↑**

| Methods | Avg. | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship |
|---|---|---|---|---|---|---|---|---|---|
| SRN [50] | 22.26 | 26.96 | 17.18 | 20.73 | 26.81 | 20.85 | 18.09 | 26.85 | 20.60 |
| NV [32] | 26.05 | 28.33 | 22.58 | 24.79 | 30.71 | 26.08 | 24.22 | 27.78 | 23.93 |
| NeRF [37] | 31.01 | 33.00 | 25.01 | 30.13 | 36.18 | 32.54 | 29.62 | 32.91 | 28.65 |
| JaxNeRF [9] | 31.69 | 34.08 | 25.03 | 30.43 | 36.92 | 33.28 | 29.91 | 34.53 | 29.36 |
| JaxNeRF+ [9] | 33.00 | 35.35 | 25.65 | 32.77 | 37.58 | 35.35 | 30.29 | 36.52 | 30.48 |
| Mip-NeRF [1] | 33.09 | 35.14 | 25.48 | 33.29 | 37.48 | 35.70 | 30.71 | 36.51 | 30.41 |
| AutoInt [29] | 25.55 | 25.60 | 20.78 | 22.47 | 32.33 | 25.09 | 25.90 | 28.10 | 24.15 |
| FastNeRF [15] | 29.97 | 32.32 | 23.75 | 27.79 | 34.72 | 32.28 | 28.89 | 31.77 | 27.69 |
| SNeRG [18] | 30.38 | 33.24 | 24.57 | 29.32 | 34.33 | 33.82 | 27.21 | 32.60 | 27.97 |
| KiloNeRF [44] | 31.00 | - | - | - | - | - | - | - | - |
| PlenOctrees [67] | 31.71 | 34.66 | 25.31 | 30.79 | 36.79 | 32.95 | 29.76 | 33.97 | 29.42 |
| NSVF [30] | 31.75 | 33.19 | 25.18 | 31.23 | 37.14 | 32.29 | 32.68 | 34.27 | 27.93 |
| Ours | 31.95 | 34.09 | 25.44 | 32.78 | 36.74 | 34.64 | 29.57 | 33.20 | 29.13 |

**SSIM↑**

| Methods | Avg. | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship |
|---|---|---|---|---|---|---|---|---|---|
| SRN [50] | 0.846 | 0.910 | 0.766 | 0.849 | 0.923 | 0.809 | 0.808 | 0.947 | 0.757 |
| NV [32] | 0.893 | 0.916 | 0.873 | 0.910 | 0.944 | 0.880 | 0.888 | 0.946 | 0.784 |
| NeRF [37] | 0.947 | 0.967 | 0.925 | 0.964 | 0.974 | 0.961 | 0.949 | 0.980 | 0.856 |
| JaxNeRF [9] | 0.953 | 0.975 | 0.925 | 0.967 | 0.979 | 0.968 | 0.952 | 0.987 | 0.868 |
| JaxNeRF+ [9] | 0.962 | 0.982 | 0.936 | 0.980 | 0.983 | 0.979 | 0.956 | 0.991 | 0.887 |
| Mip-NeRF [1] | 0.961 | 0.981 | 0.932 | 0.980 | 0.982 | 0.978 | 0.959 | 0.991 | 0.882 |
| AutoInt [29] | 0.911 | 0.928 | 0.861 | 0.898 | 0.974 | 0.900 | 0.930 | 0.948 | 0.852 |
| FastNeRF [15] | 0.941 | 0.966 | 0.913 | 0.954 | 0.973 | 0.964 | 0.947 | 0.977 | 0.805 |
| SNeRG [18] | 0.950 | 0.975 | 0.929 | 0.967 | 0.971 | 0.973 | 0.938 | 0.982 | 0.865 |
| KiloNeRF [44] | 0.95 | - | - | - | - | - | - | - | - |
| PlenOctrees [67] | 0.958 | 0.981 | 0.933 | 0.970 | 0.982 | 0.971 | 0.955 | 0.987 | 0.884 |
| NSVF [30] | 0.953 | 0.968 | 0.931 | 0.973 | 0.980 | 0.960 | 0.973 | 0.987 | 0.854 |
| Ours | 0.957 | 0.977 | 0.930 | 0.978 | 0.980 | 0.976 | 0.951 | 0.983 | 0.879 |

**LPIPS↓ (Vgg)**

| Methods | Avg. | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship |
|---|---|---|---|---|---|---|---|---|---|
| SRN [50] | 0.170 | 0.106 | 0.267 | 0.149 | 0.100 | 0.200 | 0.174 | 0.063 | 0.299 |
| NV [32] | 0.160 | 0.109 | 0.214 | 0.162 | 0.109 | 0.175 | 0.130 | 0.107 | 0.276 |
| NeRF [37] | 0.081 | 0.046 | 0.091 | 0.044 | 0.121 | 0.050 | 0.063 | 0.028 | 0.206 |
| JaxNeRF [9] | 0.068 | 0.035 | 0.085 | 0.038 | 0.079 | 0.040 | 0.060 | 0.019 | 0.185 |
| Mip-NeRF [1] | 0.043 | 0.021 | 0.065 | 0.020 | 0.027 | 0.021 | 0.040 | 0.009 | 0.138 |
| PlenOctrees [67] | 0.053 | 0.022 | 0.076 | 0.038 | 0.032 | 0.034 | 0.059 | 0.017 | 0.144 |
| Ours | 0.053 | 0.027 | 0.077 | 0.024 | 0.034 | 0.028 | 0.058 | 0.017 | 0.161 |

**LPIPS↓ (Alex)**

| Methods | Avg. | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship |
|---|---|---|---|---|---|---|---|---|---|
| NSVF [30] | 0.047 | 0.043 | 0.069 | 0.017 | 0.025 | 0.029 | 0.021 | 0.010 | 0.162 |
| Ours | 0.035 | 0.016 | 0.061 | 0.015 | 0.017 | 0.014 | 0.026 | 0.014 | 0.118 |

Table L1. Quantitative results on each scene from the **Synthetic-NeRF** [37] dataset. We highlight the top 3 results of each column under each metric in gold, silver, and bronze.
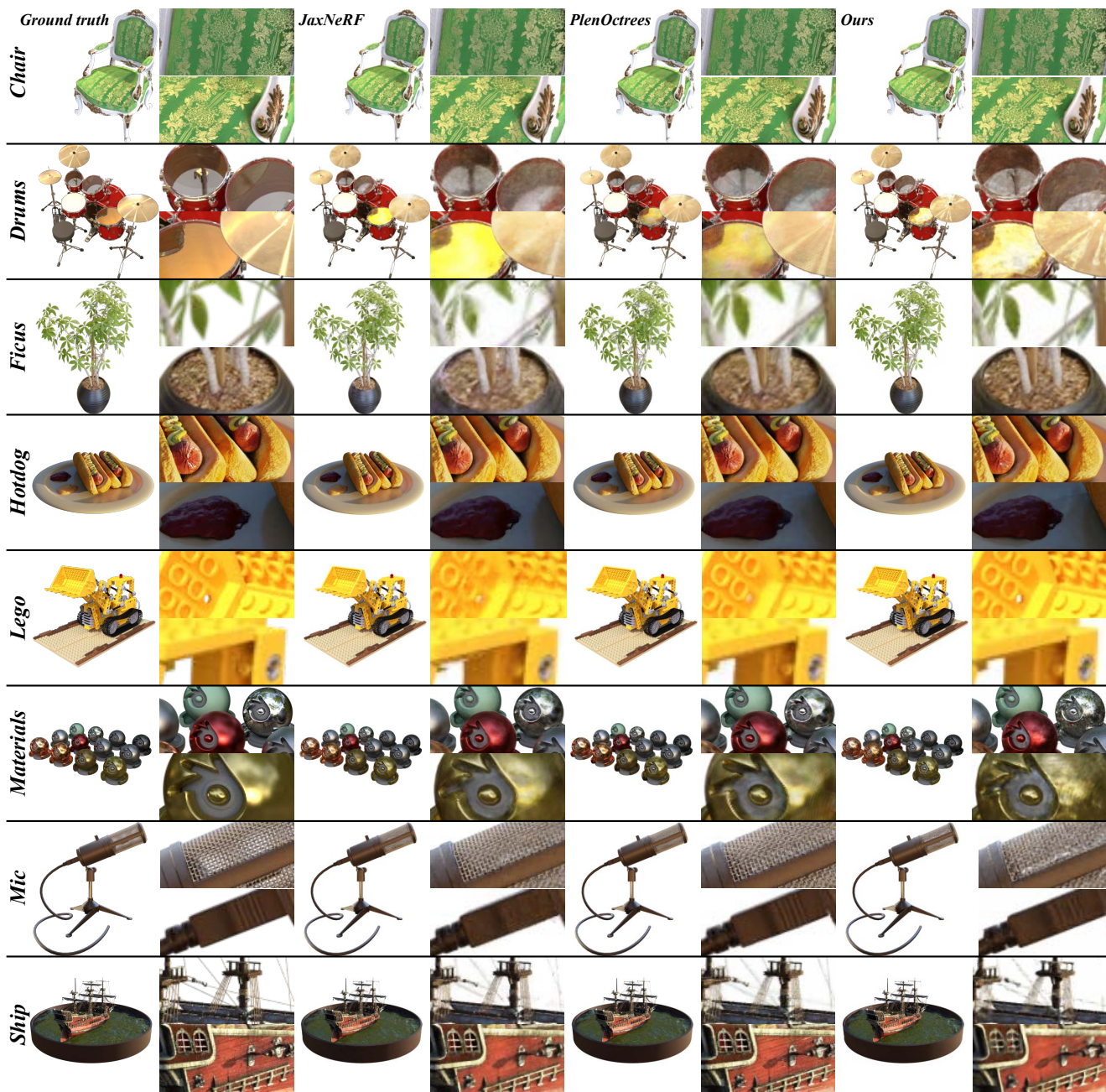
Figure L1. Qualitative comparisons on the on the **Synthetic-NeRF** [37] dataset. We manually resize, crop, and compress the images. JaxNeRF and PlenOctrees are both stronger baselines than NeRF. No method produces consistently better fine detail on all scenes. JaxNeRF recovers the shadow better in *Hotdog* and no blocking artifact in *Materials*; PlenOctrees shows better fine texture on *Mic* and *Ship*; our method reconstructs the fine detail of *Ficus* and *Lego* better.

| Methods | Avg. | *Wineholder* | *Steamtrain* | *Toad* | *Robot* | *Bike* | *Palace* | *Spaceship* | *Lifestyle* |
|---|---|---|---|---|---|---|---|---|---|
| **PSNR↑** | | | | | | | | | |
| SRN [50] | 24.33 | 20.74 | 25.49 | 25.36 | 22.27 | 23.76 | 24.45 | 27.99 | 24.58 |
| NV [32] | 25.83 | 21.32 | 25.31 | 24.63 | 24.74 | 26.65 | 26.38 | 29.90 | 27.68 |
| NeRF [37] | 30.81 | 28.23 | 30.84 | 29.42 | 28.69 | 31.77 | 31.76 | 34.66 | 31.08 |
| KiloNeRF [44] | 33.37 | - | - | - | - | - | - | - | - |
| NSVF [30] | 35.13 | 32.04 | 35.13 | 33.25 | 35.24 | 37.75 | 34.05 | 39.00 | 34.60 |
| Ours | 35.08 | 30.26 | 36.56 | 33.10 | 36.36 | 38.33 | 34.49 | 37.71 | 33.79 |
| **SSIM↑** | | | | | | | | | |
| SRN [50] | 0.882 | 0.850 | 0.923 | 0.822 | 0.904 | 0.926 | 0.792 | 0.945 | 0.892 |
| NV [32] | 0.892 | 0.828 | 0.900 | 0.813 | 0.927 | 0.943 | 0.826 | 0.956 | 0.941 |
| NeRF [37] | 0.952 | 0.920 | 0.966 | 0.920 | 0.960 | 0.970 | 0.950 | 0.980 | 0.946 |
| KiloNeRF [44] | 0.97 | - | - | - | - | - | - | - | - |
| NSVF [30] | 0.979 | 0.965 | 0.986 | 0.968 | 0.988 | 0.991 | 0.969 | 0.991 | 0.971 |
| Ours | 0.975 | 0.949 | 0.989 | 0.966 | 0.992 | 0.991 | 0.962 | 0.988 | 0.965 |
| **LPIPS↓ (Vgg)** | | | | | | | | | |
| Ours | 0.033 | 0.055 | 0.019 | 0.047 | 0.013 | 0.011 | 0.043 | 0.019 | 0.054 |
| **LPIPS↓ (Alex)** | | | | | | | | | |
| SRN [50] | 0.141 | 0.224 | 0.082 | 0.204 | 0.120 | 0.075 | 0.240 | 0.061 | 0.120 |
| NV [32] | 0.124 | 0.204 | 0.121 | 0.192 | 0.096 | 0.067 | 0.173 | 0.056 | 0.088 |
| NeRF [37] | 0.043 | 0.096 | 0.031 | 0.069 | 0.038 | 0.019 | 0.031 | 0.016 | 0.047 |
| NSVF [30] | 0.015 | 0.020 | 0.010 | 0.032 | 0.007 | 0.004 | 0.018 | 0.006 | 0.020 |
| Ours | 0.019 | 0.038 | 0.010 | 0.030 | 0.005 | 0.004 | 0.027 | 0.009 | 0.027 |

Table L2. Quantitative results on each scene from the **Synthetic-NSVF** [30] dataset. We highlight the top 3 results of each column under each metric in gold , silver , and bronze .
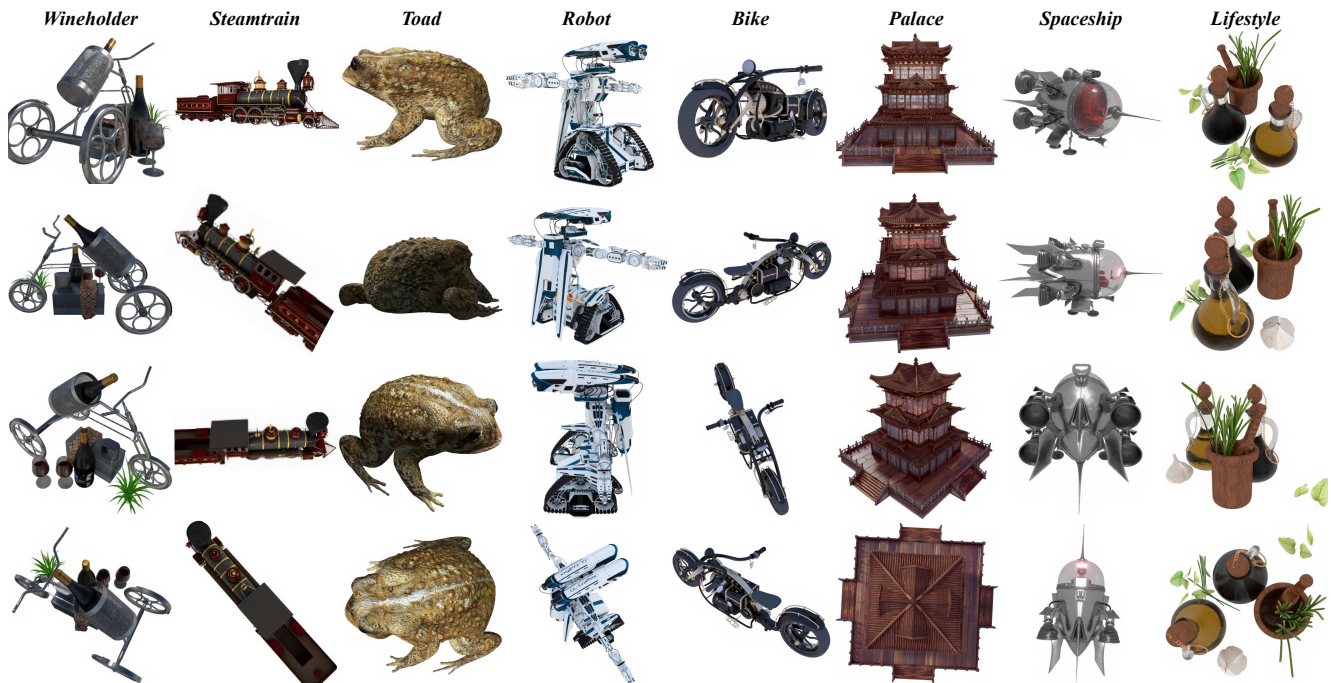


Figure L2. The synthesized novel views by our method on the **Synthetic-NSVF** [30] dataset. We manually resize, crop, and compress the images.

15

| Methods | Avg. | *Jade* | *Fountain* | *Character* | *Statues* |
|---|---|---|---|---|---|
| **PSNR↑** | | | | | |
| SRN [50] | 20.51 | 18.57 | 21.04 | 21.98 | 20.46 |
| NV [32] | 23.03 | 22.08 | 22.71 | 24.10 | 23.22 |
| NeRF [37] | 24.15 | 21.65 | 25.59 | 25.87 | 23.48 |
| KiloNeRF [44] | 27.39 | - | - | - | - |
| NSVF [30] | 26.89 | 26.96 | 27.73 | 27.95 | 24.97 |
| Ours | 28.02 | 27.68 | 28.48 | 30.31 | 25.62 |
| **SSIM↑** | | | | | |
| SRN [50] | 0.770 | 0.715 | 0.717 | 0.853 | 0.794 |
| NV [32] | 0.793 | 0.750 | 0.762 | 0.876 | 0.785 |
| NeRF [37] | 0.828 | 0.750 | 0.860 | 0.900 | 0.800 |
| KiloNeRF [44] | 0.92 | - | - | - | |
| NSVF [30] | 0.898 | 0.901 | 0.913 | 0.921 | 0.858 |
| Ours | 0.922 | 0.914 | 0.926 | 0.963 | 0.883 |
| **LPIPS↓ (Vgg)** | | | | | |
| Ours | 0.101 | 0.108 | 0.115 | 0.045 | 0.137 |
| **LPIPS↓ (Alex)** | | | | | |
| SRN [50] | 0.294 | 0.323 | 0.291 | 0.208 | 0.354 |
| NV [32] | 0.243 | 0.292 | 0.263 | 0.140 | 0.277 |
| NeRF [37] | 0.192 | 0.264 | 0.149 | 0.149 | 0.206 |
| NSVF [30] | 0.114 | 0.094 | 0.113 | 0.074 | 0.171 |
| Ours | 0.075 | 0.075 | 0.086 | 0.029 | 0.110 |

Table L3. Quantitative results on each scene from the **BlendedMVS** [66] dataset. We highlight the top 3 results of each column under each metric in gold, silver, and bronze.



Figure L3. The synthesized novel views by our method on the **BlendedMVS** [66] dataset. We manually resize, crop, and compress the images.

| Methods | Avg. | Ignatius | Truck | Barn | Caterpillar | Family |
|---|---|---|---|---|---|---|
| **PSNR↑** | | | | | | |
| SRN [50] | 24.10 | 26.70 | 22.62 | 22.44 | 21.14 | 27.57 |
| NV [32] | 23.70 | 26.54 | 21.71 | 20.82 | 20.71 | 28.72 |
| NeRF [37] | 25.78 | 25.43 | 25.36 | 24.05 | 23.75 | 30.29 |
| JaxNeRF [9] | 27.94 | 27.95 | 26.66 | 27.39 | 25.24 | 32.47 |
| KiloNeRF [44] | 28.41 | - | - | - | - | - |
| PlenOctrees [67] | 27.99 | 28.19 | 26.83 | 26.80 | 25.29 | 32.85 |
| NSVF [30] | 28.48 | 27.91 | 26.92 | 27.16 | 26.44 | 33.58 |
| Ours | 28.41 | 28.16 | 27.15 | 27.01 | 26.00 | 33.75 |
| **SSIM↑** | | | | | | |
| SRN [50] | 0.847 | 0.920 | 0.832 | 0.741 | 0.834 | 0.908 |
| NV [32] | 0.834 | 0.922 | 0.793 | 0.721 | 0.819 | 0.916 |
| NeRF [37] | 0.864 | 0.920 | 0.860 | 0.750 | 0.860 | 0.932 |
| JaxNeRF [9] | 0.904 | 0.940 | 0.896 | 0.842 | 0.892 | 0.951 |
| KiloNeRF [44] | 0.91 | - | - | - | - | - |
| PlenOctrees [67] | 0.917 | 0.948 | 0.914 | 0.856 | 0.907 | 0.962 |
| NSVF [30] | 0.901 | 0.930 | 0.895 | 0.823 | 0.900 | 0.954 |
| Ours | 0.911 | 0.944 | 0.906 | 0.838 | 0.906 | 0.962 |
| **LPIPS↓ (Vgg)** | | | | | | |
| JaxNeRF [9] | 0.168 | 0.102 | 0.173 | 0.286 | 0.189 | 0.092 |
| PlenOctrees [67] | 0.131 | 0.080 | 0.130 | 0.226 | 0.148 | 0.069 |
| Ours | 0.155 | 0.083 | 0.160 | 0.294 | 0.167 | 0.070 |
| **LPIPS↓ (Alex)** | | | | | | |
| SRN [50] | 0.251 | 0.128 | 0.266 | 0.448 | 0.278 | 0.134 |
| NV [32] | 0.260 | 0.117 | 0.312 | 0.479 | 0.280 | 0.111 |
| NeRF [37] | 0.198 | 0.111 | 0.192 | 0.395 | 0.196 | 0.098 |
| NSVF [30] | 0.155 | 0.106 | 0.148 | 0.307 | 0.141 | 0.063 |
| Ours | 0.148 | 0.090 | 0.145 | 0.290 | 0.152 | 0.064 |

Table L4. Quantitative results on each scene from the **Tanks&Temples** [21] dataset dataset. We highlight the top 3 results of each column under each metric in gold , silver , and bronze .
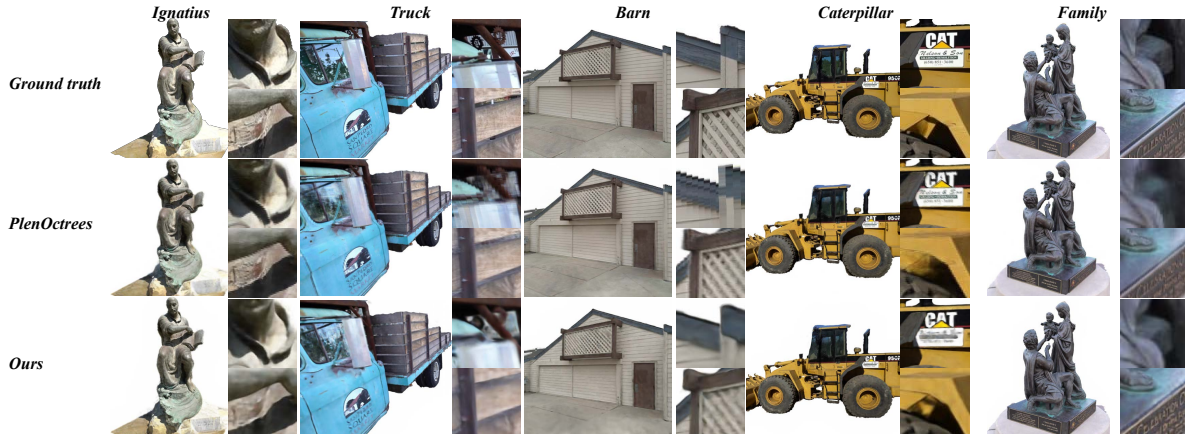


Figure L4. Qualitative comparisons on the **Tanks&Temples** [21] dataset. We manually resize, crop, and compress the images. Our quality is comparable to PlenOctrees, and no method produces consistent finer detail. Besides, we do not show blocking artifacts.

| Methods | Avg. | *Chair* | *Pedestal* | *Cube* | *Vase* |
|---|---|---|---|---|---|
| **PSNR↑** | | | | | |
| Nearest Neighbor | 20.94 | 20.69 | 21.49 | 18.32 | 23.26 |
| NV [32] | 29.62 | 35.15 | 36.47 | 26.48 | 20.39 |
| DeepVoxels [49] | 30.55 | 33.45 | 32.35 | 28.42 | 27.99 |
| DeepVoxels++ [17] | 37.31 | 40.87 | 38.93 | 36.51 | 32.91 |
| SRN [50] | 33.20 | 36.67 | 35.91 | 28.74 | 31.46 |
| LLFF [36] | 34.38 | 36.11 | 35.87 | 32.58 | 32.97 |
| NeRF [37] | 40.15 | 42.65 | 41.44 | 39.19 | 37.32 |
| IBRNet [60] | 42.93 | - | - | - | - |
| Ours | 45.83 | 48.48 | 48.51 | 43.77 | 42.54 |
| **SSIM↑** | | | | | |
| Nearest Neighbor | 0.89 | 0.94 | 0.87 | 0.83 | 0.92 |
| NV [32] | 0.929 | 0.980 | 0.963 | 0.916 | 0.857 |
| DeepVoxels [49] | 0.97 | 0.99 | 0.97 | 0.97 | 0.96 |
| DeepVoxels++ [17] | 0.99 | 0.99 | 0.98 | 0.99 | 0.98 |
| SRN [50] | 0.963 | 0.982 | 0.957 | 0.944 | 0.969 |
| LLFF [36] | 0.985 | 0.992 | 0.983 | 0.983 | 0.983 |
| NeRF [37] | 0.991 | 0.991 | 0.986 | 0.996 | 0.992 |
| IBRNet [60] | 0.997 | - | - | - | - |
| Ours | 0.998 | 0.998 | 0.998 | 0.998 | 0.998 |
| **LPIPS↓ (Vgg)** | | | | | |
| SRN [50] | 0.073 | 0.093 | 0.081 | 0.074 | 0.044 |
| NV [32] | 0.099 | 0.096 | 0.069 | 0.113 | 0.117 |
| LLFF [36] | 0.048 | 0.051 | 0.039 | 0.064 | 0.039 |
| NeRF [37] | 0.023 | 0.047 | 0.024 | 0.006 | 0.017 |
| IBRNet [60] | 0.009 | - | - | - | - |
| Ours | 0.006 | 0.015 | 0.003 | 0.002 | 0.004 |
| **LPIPS↓ (Alex)** | | | | | |
| Ours | 0.002 | 0.005 | 0.001 | 0.001 | 0.002 |

Table L5. Quantitative results on each scene from the **DeepVoxels** [49] dataset. We highlight the top 3 results of each column under each metric in gold , silver , and bronze .
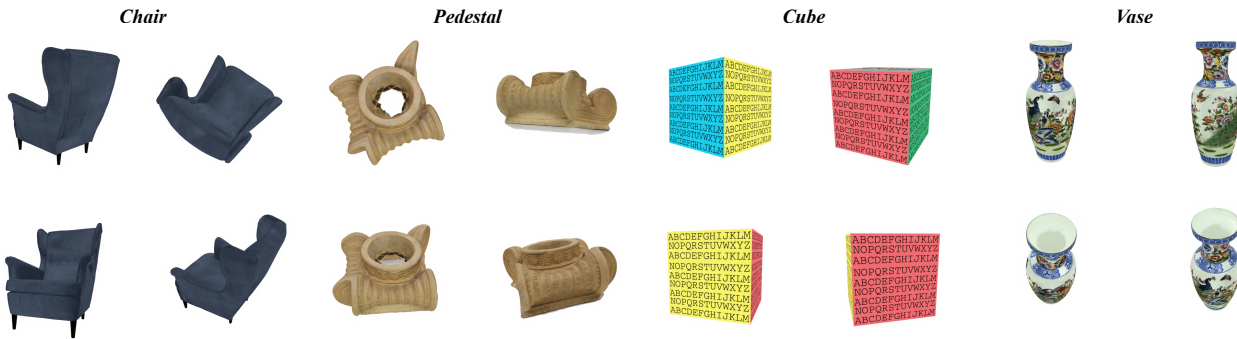


Figure L5. The synthesized novel views by our method on the **DeepVoxels** [49] dataset.

## M. Derivation of low-density initialization

In this section, we derive the bias term $b$ in the low-density initialization to make the activated alpha value very close to zero (i.e., $\alpha^{(\text{init})(\text{c})} \approx 0$) at the beginning of our coarse stage optimization. More specifically, $\alpha^{(\text{init})(\text{c})}$ is a hyperparameter, and $\left(1 - \alpha^{(\text{init})(\text{c})}\right)$ means the decay factor of the accumulated transmittance for a ray tracing forward a distance of the voxel size $s^{(\text{c})}$.

Below we omit the stage-specific superscript for simplicity. Let the accumulated distance of $s$ be divided into $N$ segments, each with a distance of $\delta_i$, i.e.,

$$\sum_{i=1}^{N} \delta_i = s .$$

Recall that, in practice, we initialize all raw values $\ddot{\sigma}$ in $V^{(\text{density})(\text{c})}$ to 0 and thus the initial activated densities $\sigma_i$ are all equal to $\log\left(1 + \exp(b)\right)$. Then, the decay factor of the accumulated transmittance can be written as

$$\prod_{i=1}^{N}(1 - \alpha_i) = \prod_{i=1}^{N}\left((1 - (1 - \exp(-\sigma_i \delta_i)))\right)$$
$$= \prod_{i=1}^{N} \exp(-\sigma_i \delta_i)$$
$$= \prod_{i=1}^{N} \exp\left(\log\left(1 + \exp(b)\right) \cdot (-\delta_i)\right)$$
$$= \exp\left(\sum_{i=1}^{N}\left(\log\left(1 + \exp(b)\right) \cdot (-\delta_i)\right)\right)$$
$$= \exp\left(\log\left(1 + \exp(b)\right) \cdot \sum_{i=1}^{N}(-\delta_i)\right)$$
$$= \exp\left(\log\left(1 + \exp(b)\right) \cdot (-s)\right)$$
$$= (1 + \exp(b))^{-s} .$$

We want to find the $b$ such that the decay factor is $(1 - \alpha^{(\text{init})})$ for passing through a distance of the voxel size $s$. Therefore, we have

$$b = \log\left(\left(1 - \alpha^{(\text{init})}\right)^{-\frac{1}{s}} - 1\right) ,$$

which is the equation presented in our main paper to impose the prior of low-density initialization.

## N. Derivation of post-activation

By expanding the post-activated trilinear interpolation, we have

$$\alpha^{(\text{post})} = 1 - \exp\left(-\delta \log\left(1 + \exp\left(\text{interp}\left(\boldsymbol{x}, \boldsymbol{V}\right)\right)\right)\right)$$
$$= 1 - \left(1 + \exp\left(\text{interp}\left(\boldsymbol{x}, \boldsymbol{V}\right)\right)\right)^{-\delta} ,$$

$$(15)$$

where $\delta$ is a volume rendering related value and is treated as a constant in later derivation, and $\text{interp}$ is the interpolation operation of a spatial point $\boldsymbol{x}$ in grid $\boldsymbol{V}$. The bias term $b$ is omitted here for simplicity. The overall activation maps an interpolant ($\text{interp}\left(\boldsymbol{x}, \boldsymbol{V}\right) \in \mathbb{R}$) into an alpha value ($\alpha^{(\text{post})} \in [0, 1]$).

We want to prove that such a post-activation can produce sharp linear surface (decision boundary) in a single grid cell. Let first prove in the simplest 1D grid and then we generalize it to 2D grid. The case in 3D grid then can be proven easily using the derivations in 1D and 2D grid.
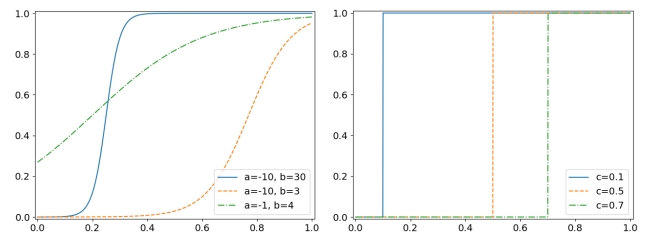
### N.1. Derivation for a 1D grid cell

In 1D grid, $\boldsymbol{x}$ is a scalar. Without loss of generality, we assume $\boldsymbol{x} = 0$ and $\boldsymbol{x} = 1$ are the left and right bound of the 1D grid cell, respectively. Let $a, b \in \mathbb{R}$ be the grid values stored at the cell's left and right bound. Then Eq. (15) with 1D linear interpolation of the cell can be re-written for this specific case:

$$S(\boldsymbol{x}; a, b) = 1 - (1 + \exp(a(1 - \boldsymbol{x}) + b\boldsymbol{x}))^{-\delta} . \quad (16)$$

Consider a target function $T(\boldsymbol{x}; c)$ in the form of a shifted unit step function,

$$T(\boldsymbol{x}; c) = \mathbb{1}(\boldsymbol{x} - c)$$
$$= \begin{cases} 1, & \boldsymbol{x} > c, \\ 0, & \boldsymbol{x} \le c, \end{cases} \quad (17)$$

where $0 < c < 1$ is the position of the target linear surface (decision boundary) in the 1D grid cell. We only derive for the occupancy at right-hand side as the opposite direction can be trivially generalized. Visualization for $S$ and $T$ with some specific parameters is shown in Fig. N1.



(a) Three examples of $S(\boldsymbol{x}; a, b)$ with $\delta = 1$

(b) Three examples of $T(\boldsymbol{x}; c)$

Figure N1. Visualization for some examples of $S$ and $T$ with different parameter settings.

We show that $S(\boldsymbol{x}; a, b)$ can be made arbitrarily close to $T(\boldsymbol{x}; c)$. Specifically, given any $\epsilon$ and $\Delta$ satisfying $0 < \epsilon < 1$ and $0 < \Delta < \min(c, 1 - c)$, we can find the grid values $a, b$ such that

$$|S(\boldsymbol{x}; a, b) - T(\boldsymbol{x}; c)| \le \epsilon, \ \forall |\boldsymbol{x} - c| \ge \Delta .$$

As the function $S$ is a monotonically increasing function bounded in $[0, 1]$, the criterion can then be simplified into

$$1 - S(c + \Delta;\ a, b) \leq \epsilon\ , \tag{18a}$$
$$S(c - \Delta;\ a, b) \leq \epsilon\ . \tag{18b}$$

By expanding the inequality (18a), we get

$$1 - S(c + \Delta;\ a, b)$$
$$= 1 - \left(1 - (1 + \exp(a(1 - (c + \Delta)) + b(c + \Delta))^{-\delta}\right)$$
$$= (1 + \exp(a(1 - (c + \Delta)) + b(c + \Delta))^{-\delta} \leq \epsilon\ .$$

We solve the inequality:

$$\left(\frac{1}{1 + \exp(a(1 - (c + \Delta)) + b(c + \Delta))}\right)^{\delta} \leq \epsilon$$
$$\frac{1}{1 + \exp(a(1 - (c + \Delta)) + b(c + \Delta))} \leq \epsilon^{1/\delta}$$
$$1 + \exp(a(1 - (c + \Delta)) + b(c + \Delta)) \geq \frac{1}{\epsilon^{1/\delta}}$$
$$\exp(a(1 - (c + \Delta)) + b(c + \Delta)) \geq \frac{1}{\epsilon^{1/\delta}} - 1\ .$$

Finally, we obtain

$$a(1 - (c + \Delta)) + b(c + \Delta) \geq \log\left(\frac{1}{\epsilon^{1/\delta}} - 1\right)\ . \tag{19}$$

By applying the same process to the inequality (18b), we get another inequality:

$$a(1 - (c - \Delta)) + b(c - \Delta) \leq \log\left(\frac{1}{(1 - \epsilon)^{1/\delta}} - 1\right)\ . \tag{20}$$

There are infinite numbers of solutions satisfying the inequalities (19) and (20). Here, we introduce one more constraint to make the derivation and later extensions simpler:

$$S(c;\ a, b) = 0.5\ . \tag{21}$$

From Eq. (16) with this constraint we can then derive the linear relation between $a, b$:

$$b = a\frac{c - 1}{c} + \frac{\log(2^{\frac{1}{\delta}} - 1)}{c}\ . \tag{22}$$

Substitute $b$ in the inequality (19), we get

$$a(1 - (c + \Delta)) + b(c + \Delta)$$
$$= a(1 - (c + \Delta)) + \left(a\frac{c - 1}{c} + \frac{\log(2^{\frac{1}{\delta}} - 1)}{c}\right)(c + \Delta)$$
$$= a\frac{-\Delta}{c} + \log(2^{\frac{1}{\delta}} - 1)\frac{c + \Delta}{c} \geq \log\left(\frac{1}{\epsilon^{1/\delta}} - 1\right)\ .$$

Similarly, substitute $b$ in the inequality (20), we get

$$a(1 - (c - \Delta)) + b(c - \Delta)$$
$$= a(1 - (c - \Delta)) + \left(a\frac{c - 1}{c} + \frac{\log(2^{\frac{1}{\delta}} - 1)}{c}\right)(c - \Delta)$$
$$= a\frac{\Delta}{c} + \log(2^{\frac{1}{\delta}} - 1)\frac{c - \Delta}{c} \leq \log\left(\frac{1}{(1 - \epsilon)^{1/\delta}} - 1\right)\ .$$

In summary, by adding the extra constraint (21), the inequality (19) and (20) become

$$\begin{cases} a \leq & \log(2^{\frac{1}{\delta}} - 1)\frac{c + \Delta}{\Delta} - \log\left(\frac{1}{\epsilon^{1/\delta}} - 1\right)\frac{c}{\Delta} \\ a \leq & \log\left(\frac{1}{(1-\epsilon)^{1/\delta}} - 1\right)\frac{c}{\Delta} - \log(2^{\frac{1}{\delta}} - 1)\frac{c - \Delta}{\Delta}\ , \end{cases} \tag{23}$$

which defines the upper bound of $a$ such that the conditions on the function $S$ in (18a), (18b), and (21) can all be satisfied. The tighter upper bound is

$$a^{(\text{upper bound})} =$$
$$\begin{cases} \log(2^{\frac{1}{\delta}} - 1)\frac{c + \Delta}{\Delta} - \log\left(\frac{1}{\epsilon^{1/\delta}} - 1\right)\frac{c}{\Delta}, & \text{if } \delta < 1 \\ \log\left(\frac{1}{(1-\epsilon)^{1/\delta}} - 1\right)\frac{c}{\Delta} - \log(2^{\frac{1}{\delta}} - 1)\frac{c - \Delta}{\Delta}, & \text{otherwise} \end{cases}, \tag{24}$$

where the $\delta > 0$ is the pre-defined step size in volume rendering and we skip the detailed derivation of Eq. (24) here.

As a verification, we re-use the examples in Fig. N1b as the target functions and set the tolerance to $\epsilon = 10^{-4}, \Delta = 10^{-2}$ and the volume rendering related value to $\delta = 0.5$. We can directly find the grid values $a, b$ using the derivations given above. We show the derived numbers and the resulting plot visualization in Fig. N2. It can be seen that the derived $S(\boldsymbol{x};\ a, b)$ can faithfully resemble the target $T(\boldsymbol{x};\ c)$.

| $T(\boldsymbol{x};\ c)$ | By Eq. (24) | By Eq. (22) |
| --- | --- | --- |
| $c = 0.1$ | $a \approx -172.1$ | $b \approx 1560.1$ |
| $c = 0.5$ | $a \approx -865.0$ | $b \approx 867.2$ |
| $c = 0.7$ | $a \approx -1211.4$ | $b \approx 520.8$ |



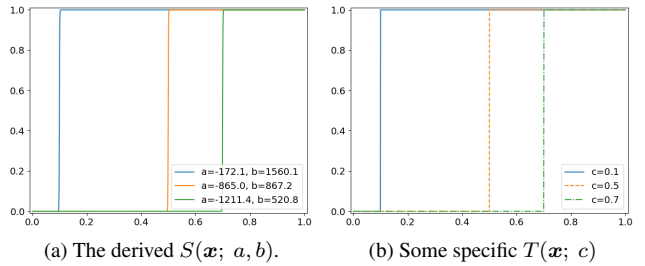(a) The derived $S(\boldsymbol{x};\ a, b)$.   (b) Some specific $T(\boldsymbol{x};\ c)$

Figure N2. Using the derived upper bound (24) and Eq. (22) to directly find the grid values $a, b$ that fit the target function $T(\boldsymbol{x};\ c)$.
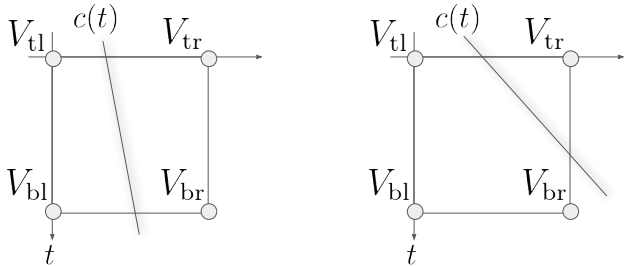
## N.2. Derivation for a 2D grid cell

We illustrate two situations that a linear boundary crossing a 2D grid cell in Fig. N3, where $V_{tl}, V_{tr}, V_{bl}, V_{br}$ are the top-left, top-right, bottom-left, and bottom-right values of a 2D grid cell and $c(t)$ is the linear boundary parameterized by the vertical position $t$. Let the coordinates of the top-left, top-right, bottom-left, and bottom-right corners be $(0,0)$, $(1,0)$, $(0,1)$, and $(1,1)$, so the top and bottom edges of the grid cell are on the horizontal lines $t = 0$ and $t = 1$, respectively. Without loss of generality, we assume the linear boundary always intersects the top edge of the 2D grid cell (*i.e.*, $0 < c(0) < 1$) and the left-hand-side $[0, c(0))$ of the boundary is free space. We can generalize to other cases via rotation and flipping, or by negating the grid values. We consider the 2D target function $T$ with respect to the decision boundary inside the grid cell as a 2D unit step function, and parameterize it by the vertical coordinate $t$ so that on each horizontal scan line the target function can be expressed as

$$T(\boldsymbol{x}(t);\ c(t)) = \mathbb{1}(\boldsymbol{x}(t) - c(t))$$
$$= \begin{cases} 1, & \boldsymbol{x}(t) > c(t), \\ 0, & \boldsymbol{x}(t) \leq c(t). \end{cases} \quad (25)$$

The goal here is to show that, by choosing suitable values for $V_{tl}, V_{tr}, V_{bl}, V_{br}$, we are able to approximate the target function closely enough within the required tolerance using our post-activation scheme.



(a) **Case I:** The boundary intersects the top and bottom edge of the grid.

(b) **Case II:** The boundary intersects the top and right edge of the grid.

Figure N3. Two cases of a linear boundary $c(t)$ crossing a 2D grid cell, where $V_{lt}, V_{rt}, V_{lb}, V_{rb}$ are the grid values of the cell. Without loss of generality, we assume the linear boundary always intersects the top edge of the grid cell, *i.e.*, $0 < c(0) < 1$.

**Case I:** $0 < c(0) < 1$ **and** $0 < c(1) < 1$. An example is illustrated in Fig. N3a, where the linear boundary intersects the top edge and the bottom edge of a grid cell. The linear boundary $c$ is a linear function of $t$ defined by $c(t) = (1 - t) \cdot c(0) + t \cdot c(1)$. Based on the results we have derived for 1D, to prove this 2D case we only need to show that Eq. (22) and Eq. (24) are linear function of $c$, so that once we use the two equations to determine $V_{tl}, V_{tr}$ for approximating

$c(0)$ and $V_{bl}, V_{br}$ for approximating $c(1)$, we can readily recover $c(t)$ from $c(0)$ and $c(1)$ on all horizontal scan lines $0 \leq t \leq 1$. That is, the approximation criteria in Eq. (18a) and Eq. (18b), where the target surface position is $c = c(t)$, are automatically satisfied by

$$a = a(t) = V_{tl}(1 - t) + V_{bl}t ,$$
$$b = b(t) = V_{tr}(1 - t) + V_{br}t .$$

Eq. (24) is trivially a linear function of $c$ given $\delta$. By substituting $a$ in Eq. (22) with $a^{(\text{upper bound})}$ from Eq. (24), we get

$$b = \log(2^{\frac{1}{\delta}} - 1)\frac{c + \Delta - 1}{\Delta} - \log\left(\frac{1}{\epsilon^{1/\delta}} - 1\right)\frac{c - 1}{\Delta}$$

when $0 < \delta < 1$, and

$$b = \log\left(\frac{1}{(1 - \epsilon)^{1/\delta}} - 1\right)\frac{c - 1}{\Delta} - \log(2^{\frac{1}{\delta}} - 1)\frac{c - \Delta - 1}{\Delta}$$

when $1 \leq \delta$. Thus, $b$ is also a linear function of $c$ given $\delta$.

**Case II:** $0 < c(0) < 1$ **and** $1 < c(1)$. We assume $0 < c < 1$ in the target function in Eq. (17), but it finally turns out that the derivations in Sec. N.1 can trivially generalize to $1 < c$. Actually, the derivations also work for $c < 0$ as long as we modify the signs in inequality (23) to '$\geq$', and the upper bound in Eq. (24) becomes lower bound. We verify the results with some examples shown in Fig. N4.

| $T(\boldsymbol{x};\ c)$ | By Eq. (24) | By Eq. (22) |
|---|---|---|
| $c = -0.6$ | $a \approx 1040.4$ | $b \approx 2772.6$ |
| $c = 1.3$ | $a \approx -2250.8$ | $b \approx -518.6$ |
| $c = 1.5$ | $a \approx -2597.2$ | $b \approx -865.0$ |



(a) The derived $S(\boldsymbol{x};\ a, b)$.

(b) Some specific $T(\boldsymbol{x};\ c)$.
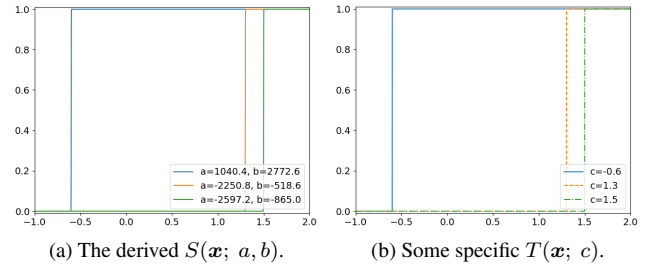
Figure N4. The derivation in Sec. N.1 is also applicable to extrapolation ($c < 0$ or $c > 1$) with minor modifications.

In summary, the derivation in Sec. N.1 also works for extrapolation with minor modifications and the derivation in **Case I** is directly applicable for **Case II**.

As a verification, we construct two examples in Fig. N5. The volume rendering step size is set to $\delta = 0.5$, the error bound is set to $\epsilon = 10^{-4}$, and we show two results with tolerance $\Delta = 0.20$ and $\Delta = 0.05$. For each target boundary, we first evaluate $c(0)$ and $c(1)$. The values of $V_{tl}, V_{bl}$ can then be derived from Eq. (24) and the values of $V_{tr}, V_{br}$ can be derived from Eq. (22).
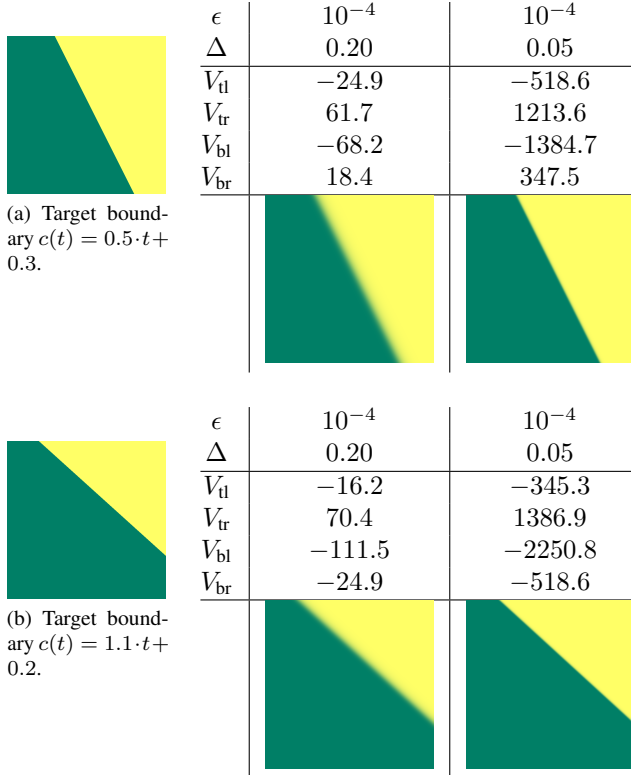
| $\epsilon$ | $10^{-4}$ | $10^{-4}$ |
|---|---|---|
| $\Delta$ | 0.20 | 0.05 |
| $V_{\mathrm{tl}}$ | $-24.9$ | $-518.6$ |
| $V_{\mathrm{tr}}$ | 61.7 | 1213.6 |
| $V_{\mathrm{bl}}$ | $-68.2$ | $-1384.7$ |
| $V_{\mathrm{br}}$ | 18.4 | 347.5 |

(a) Target boundary $c(t) = 0.5 \cdot t + 0.3$.

| $\epsilon$ | $10^{-4}$ | $10^{-4}$ |
|---|---|---|
| $\Delta$ | 0.20 | 0.05 |
| $V_{\mathrm{tl}}$ | $-16.2$ | $-345.3$ |
| $V_{\mathrm{tr}}$ | 70.4 | 1386.9 |
| $V_{\mathrm{bl}}$ | $-111.5$ | $-2250.8$ |
| $V_{\mathrm{br}}$ | $-24.9$ | $-518.6$ |

(b) Target boundary $c(t) = 1.1 \cdot t + 0.2$.

Figure N5. Using the extended 2D derivation to directly find the grid values $V_{\mathrm{tl}}, V_{\mathrm{tr}}, V_{\mathrm{bl}}, V_{\mathrm{br}}$ that fit the target boundary $c(t)$ under the required error bound $\epsilon$ and tolerance $\Delta$.

## N.3. Derivation for a 3D grid cell

Similar to the proof for 2D in Sec. N.2, we can assume that the 3D linear surface $c(t, u)$ intersects the top face of a 3D grid cell, as illustrated in Fig. N6. Without loss of generality, we assume the surface $c(t, u)$ intersects the horizontal planes $u = 0$ and $u = 1$, and the left-hand-side of the surface,

$$\{[t, u, v] \mid 0 \le t \le 1, 0 \le u \le 1, 0 \le v \le s(t, u)\} \,,$$

is free space. The linear surface is $c(t, u) = (1 - u) \cdot c(t, 0) + u \cdot c(t, 1)$. We can use the results in Sec. N.2 to determine the grid values of the top four corners with $c(t, 0)$ and the values of the bottom four corners with $c(t, 1)$. The linear boundary at every horizontal slice $0 \le u \le 1$ can be automatically satisfied, as we have shown in Sec. N.2.

## N.4. Future extensions

**Beyond linear surface (decision boundary).** We only prove that the alpha field by post-activated interpolation can be arbitrarily close to a linear surface. We show in Fig. N7 that we can further tune the tolerances at the top edge and bottom edge of a grid to produce sharp and non-linear surfaces.
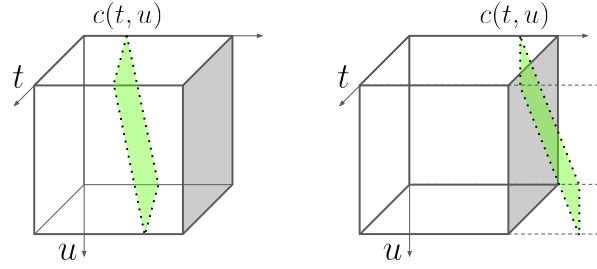


Figure N6. A linear surface $c(t, u)$ crossing a 3D grid cells. Without loss of generality, we assume the surface intersects with two horizontals planes, $u = 0$ and $u = 1$, which the top face and the bottom face of a 3D grid cell aligned with.
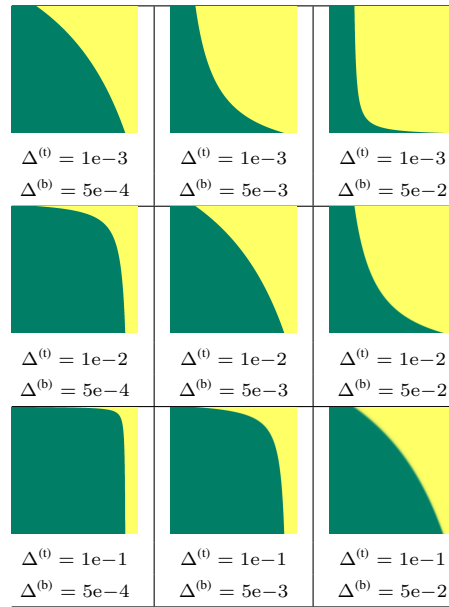


Figure N7. We fix $c(0) = 0.2, c(1) = 0.9, \delta = 0.5, \epsilon = 10^{-4}$ but use different tolerances for the top edge ($\Delta^{(\mathrm{t})}$) and the bottom edge ($\Delta^{(\mathrm{b})}$). Following the same procedure as in Sec. N.2 to determine the grid values, we can obtain a sharp non-linear surface.

Extending the proof or the capability of post-activation in future work could be helpful to geometry modeling.

**Closed form solution when 3D available.** In this work, we only consider the same input setup as NeRF, where only 2D observations and camera poses are available. In cases that the 3D model of the scene is available, an algorithm to convert the 3D model to our post-activated density voxel grid can be helpful. Our representation is directly compatible with gradient-based optimization and volume rendering to support follow-up applications, while 3D in other formats like mesh or point cloud may need more effort. We believe our derivations are useful for future work to develop a closed-form solution to convert the 3D in other formats into our representation.

# References

[1] 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, 3, 4, 7, 9, 11, 12, 13

[2] Sai Bi, Zexiang Xu, Pratul P. Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Milos Hasan, Yannick Hold-Geoffroy, David J. Kriegman, and Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition. *arxiv CS.CV 2106.01970*, 2020. 2

[3] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P. A. Lensch. Nerd: Neural reflectance decomposition from image collections. In *ICCV*, 2021. 2

[4] Chris Buehler, Michael Bosse, Leonard McMillan, Steven J. Gortler, and Michael F. Cohen. Unstructured lumigraph rendering. In *SIGGRAPH*, 2001. 2

[5] Eric R. Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. Pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021. 2

[6] 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. 1, 3, 7

[7] Abe Davis, Marc Levoy, and Frédo Durand. Unstructured light fields. *Comput. Graph. Forum*, 2012. 2

[8] Paul E. Debevec, Camillo J. Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry- and image-based approach. In *SIGGRAPH*, 1996. 2

[9] Boyang Deng, Jonathan T. Barron, and Pratul P. Srinivasan. JaxNeRF: an efficient JAX implementation of NeRF, 2020. 6, 7, 11, 13, 17

[10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. *arxiv CS.CV 2107.02791*, 2021. 1, 3

[11] Helisa Dhamo, Keisuke Tateno, Iro Laina, Nassir Navab, and Federico Tombari. Peeking behind objects: Layered depth prediction from a single image. *Pattern Recognit. Lett.*, 2019. 2

[12] John Flynn, Michael Broxton, Paul E. Debevec, Matthew DuVall, Graham Fyffe, Ryan S. Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *CVPR*, 2019. 2

[13] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, 2021. 2

[14] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, 2021. 2

[15] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien P. C. Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. 1, 2, 3, 7, 13

[16] Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. The lumigraph. In *SIGGRAPH*, 1996. 2

[17] Tong He, John P. Collomosse, Hailin Jin, and Stefano Soatto. Deepvoxels++: Enhancing the fidelity of novel view synthesis from 3d voxel embeddings. In Hiroshi Ishikawa, Cheng-Lin Liu, Tomás Pajdla, and Jianbo Shi, editors, *ACCV*, 2020. 2, 18

[18] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021. 1, 2, 3, 5, 7, 10, 13

[19] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Animashree Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *ICCV*, 2021. 2

[20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6

[21] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: benchmarking large-scale scene reconstruction. *ACM Trans. Graph.*, 2017. 6, 11, 12, 17

[22] Adam R. Kosiorek, Heiko Strathmann, Daniel Zoran, Pol Moreno, Rosalia Schneider, Sona Mokrá, and Danilo Jimenez Rezende. Nerf-vae: A geometry aware 3d scene generative model. In *ICML*, 2021. 2

[23] Anat Levin and Frédo Durand. Linear view synthesis using a dimensionality gap light field prior. In *CVPR*, 2010. 2

[24] Marc Levoy and Pat Hanrahan. Light field rendering. In *SIGGRAPH*, 1996. 2

[25] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 2

[26] Zhengqi Li, Wenqi Xian, Abe Davis, and Noah Snavely. Crowdsampling the plenoptic function. In *ECCV*, 2020. 2

[27] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. BARF: bundle-adjusting neural radiance fields. In *ICCV*, 2021. 2

[28] Yen-Chen Lin, Pete Florence, Jonathan T. Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. inerf: Inverting neural radiance fields for pose estimation. In *IROS*, 2021. 2

[29] David B. Lindell, Julien N. P. Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *CVPR*, 2021. 1, 7, 13

[30] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 1, 2, 5, 6, 7, 11, 12, 13, 15, 16, 17

[31] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. *arxiv CS.CV 2107.13421*, 2021. 1, 3, 7, 10

[32] Stephen Lombardi, Tomas Simon, Jason M. Saragih, Gabriel Schwartz, Andreas M. Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 2019. 2, 7, 13, 15, 16, 17, 18

[33] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 2

[34] Nelson L. Max. Optical models for direct volume rendering. *IEEE Trans. Vis. Comput. Graph.*, 1995. 3

[35] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *ICCV*, 2021. 2

[36] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: practical view synthesis with prescriptive sampling guidelines. *ACM Trans. Graph.*, 2019. 2, 18

[37] 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, 3, 6, 7, 9, 11, 12, 13, 14, 15, 16, 17, 18

[38] Atsuhiro Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance field. In *ICCV*, 2021. 2

[39] Michael Oechsle, Songyou Peng, and Andreas Geiger. UNISURF: unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *ICCV*, 2021. 9

[40] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. In *ICCV*, 2021. 2

[41] 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 CS.CV 2106.13228*, 2021. 2

[42] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2

[43] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. In *CVPR*, 2021. 1

[44] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, 2021. 1, 3, 7, 13, 15, 16, 17

[45] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: generative radiance fields for 3d-aware image synthesis. In *NeurIPS*, 2020. 2

[46] Jonathan Shade, Steven J. Gortler, Li-wei He, and Richard Szeliski. Layered depth images. In *SIGGRAPH*, 1998. 2

[47] Lixin Shi, Haitham Hassanieh, Abe Davis, Dina Katabi, and Frédo Durand. Light field reconstruction using sparsity in the continuous fourier domain. *ACM Trans. Graph.*, 2014. 2

[48] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *CVPR*, 2020. 2

[49] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 2, 6, 7, 11, 12, 18

[50] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 7, 13, 15, 16, 17, 18

[51] Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, 2021. 2

[52] Pratul P. Srinivasan, Richard Tucker, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *CVPR*, 2019. 2

[53] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P. Srinivasan, Jonathan T. Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *CVPR*, 2021. 2

[54] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *NeurIPS*, 2020. 3

[55] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: image synthesis using neural textures. *ACM Trans. Graph.*, 2019. 2

[56] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a deforming scene from monocular video. In *ICCV*, 2021. 2

[57] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *CVPR*, 2020. 2

[58] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3d scene inference via view synthesis. In *ECCV*, 2018. 2

[59] Michael Waechter, Nils Moehrle, and Michael Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *ECCV*, 2014. 2

[60] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas A. Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 1, 3, 7, 18

[61] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 2004. 6

[62] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural radiance fields without known camera parameters. *arxiv CS.CV 2102.07064*, 2021. 2

[63] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2, 3, 10

[64] Daniel N. Wood, Daniel I. Azuma, Ken Aldinger, Brian Curless, Tom Duchamp, David Salesin, and Werner Stuetzle. Surface light fields for 3d photography. In *SIGGRAPH*, 2000. 2

[65] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *CVPR*, 2021. 2

[66] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *CVPR*, 2020. 6, 11, 12, 16

[67] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 1, 2, 3, 5, 7, 8, 10, 11, 13, 17

[68] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. 3

[69] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arxiv CS.CV 2010.07492*, 2020. 3

[70] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 6

[71] Xiuming Zhang, Pratul P. Srinivasan, Boyang Deng, Paul E. Debevec, William T. Freeman, and Jonathan T. Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *arxiv CS.CV 2106.01970*, 2021. 2

[72] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM Trans. Graph.*, 2018. 2