

Tetra-NeRF: Representing Neural Radiance Fields Using Tetrahedra

Jonas Kulhanek
Czech Technical University in Prague
jonas.kulhanek@cvut.cz

Torsten Sattler
Czech Technical University in Prague
torsten.sattler@cvut.cz

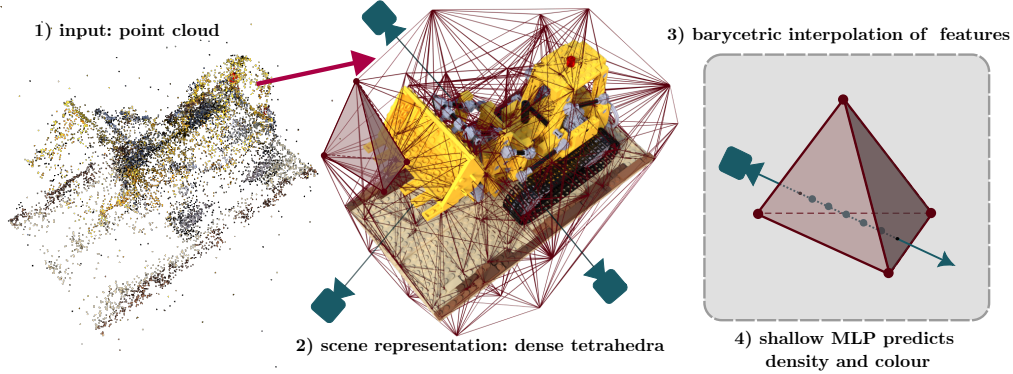


Figure 1. The input to Tetra-NeRF is a point cloud which is triangulated to get a set of tetrahedra used to represent the radiance field. Rays are sampled, and the field is queried. The barycentric interpolation is used to interpolate tetrahedra vertices, and the resulting features are passed to a shallow MLP to get the density and colours for volumetric rendering.

Abstract

Neural Radiance Fields (NeRFs) are a very recent and very popular approach for the problems of novel view synthesis and 3D reconstruction. A popular scene representation used by NeRFs is to combine a uniform, voxel-based subdivision of the scene with an MLP. Based on the observation that a (sparse) point cloud of the scene is often available, this paper proposes to use an adaptive representation based on tetrahedra obtained by the Delaunay triangulation instead of the uniform subdivision or point-based representations. We show that such a representation enables efficient training and leads to state-of-the-art results. Our approach elegantly combines concepts from 3D geometry processing, triangle-based rendering, and modern neural radiance fields. Compared to voxel-based representations, ours provides more detail around parts of the scene likely to be close to the surface. Compared to point-based representations, our approach achieves better performance.

1. Introduction

Reconstructing 3D scenes from images and rendering photo-realistic novel views is a key problem in computer vision. Recently, NeRFs [3, 4, 38] became dominant in the field for their superior photo-realistic results. Originally, NeRFs used MLPs to represent the 3D scene as an im-

PLICIT function. Given a set of posed images, NeRF randomly samples a batch of pixels, casts rays from the pixels into the 3D space, queries the implicit function at randomly sampled distances along the ray, and aggregates the sampled values using volumetric rendering [36, 38]. While the visual results of such methods are of high quality, the problem is that querying large MLPs at millions of points is costly. Also, once the network is trained, it is difficult to make any changes to the represented radiance field as everything is baked into the MLPs parameters, and any change has a non-local effect. Since then, there have been a lot of proposed alternatives to the large MLP field representation [9, 17, 31, 40, 54, 67, 71]. These methods combine an MLP with a voxel feature grid [40, 54], or in some cases represent the radiance field directly as a tensor [9, 10, 17]. When querying these representations, first, the containing voxel is found, and the features stored at the eight corner points of the voxel are trilinearly interpolated. The result is either passed through a shallow MLP [9, 10, 40, 54] or is used directly as the density and colour [17, 31, 54].

Having a dense tensor represent the entire scene is very inefficient, as we only need to represent a small space around surfaces. Therefore, different methods propose different ways of tackling the issue. Instant-NGP [40], for example, uses a hash grid instead of a dense tensor, where it relies on optimisation to resolve the hash collisions. However, similarly to MLPs, any change to the stored hashmap

influences the field in many places. A more common direction to addressing the issue is by directly using a sparse tensor representation [9, 17]. These methods start with a low-resolution grid and, at predefined steps, subsample the representation, increasing the resolution. These approaches tend to require a careful setting of hyperparameters, such as the scene bounding box and the subdivision steps, in order for the methods to work well.

Because many of these methods use traditional structure from motion (SfM) [52, 53] methods to generate the initial poses for the captured images, we can reuse the original reconstruction in the scene representation. Inspired by classical surface reconstruction methods [22, 23, 28, 29], we represent the scene as a dense triangulation of the input point cloud, where the scene is a set of non-overlapping tetrahedra whose union is the convex hull of the original point cloud [14]. When querying such a representation, we find to which tetrahedron the query point belongs and perform barycentric linear interpolation of the features stored in the vertices of the tetrahedron. This very simple representation can be thought of as the direct extension of the classical triangle-rendering pipelines used in graphics [39, 42, 43]. The representation avoids problems with the sparsity of the input point cloud as the tetrahedra fully cover the scene, resulting in a continuous rather than discrete representation.

This paper makes the following contributions: (1) We propose a novel radiance field representation which is initialised from a sparse or dense point cloud. This representation is naturally denser in the proximity of surfaces and, therefore, provides a higher resolution in these regions. (2) The proposed representation is evaluated on multiple synthetic and real-world datasets and is compared with a state-of-the-art point-cloud-based representation – Point-NeRF [67]. The presented results show that our method clearly outperforms this baseline. We further demonstrate the effectiveness of our adaptive representation by comparing it with a voxel-based representation that uses the same number of trainable parameters. (3) We make the source code and model checkpoints publicly available.¹

2. Related work

Multi-view reconstruction. The problem of multi-view reconstruction has been studied extensively and tackled with a variety of structure from motion (SfM) [52, 58, 60], and multi-view stereo (MVS) [12, 18, 53, 68] methods. These methods usually output the scene represented as a point cloud [52, 53]. In most rendering approaches, the point cloud is converted into a mesh [24, 34], and novel views are rendered by reprojecting observed images into each novel viewpoint and blending them together using either heuristically-defined [8, 13, 66] or learned [20, 48, 49,

70] blending weights.

However, the process of getting the meshes is usually quite noisy, and the resulting meshes tend to have inaccurate geometry in regions with fine details or complex materials. Instead of using noisy meshes, point-based neural rendering methods [1, 26, 37, 51] perform splatting of neural features and use 2D convolutions to render them. In contrast to these methods, our approach operates and aggregates features directly in 3D and does not suffer from the noise in the point cloud or the reconstructed mesh.

Neural radiance fields. Recently, NeRFs [3, 4, 33, 38, 74] have gained a lot of attention thanks to their high-quality rendering performance. The original NeRF method [38] was extended to better handle aliasing artefacts in [3], to better represent unbounded scenes in [4, 46, 74], or to handle real-world captured images [35, 55]. The training of the large MLPs used in these methods can be quite slow, and there has been a lot of effort on speeding up either the training [9, 17, 40] or the rendering [21, 45, 46, 71] sometimes at the cost of larger storage requirements. Other approaches tackled different aspects of NeRFs like view-dependent artefacts [59], relighting [5, 7], or proposed generative models [44, 61]. Also, a popular research direction is making the models generalize across different scenes [11, 27, 47, 63, 72]. A large area of research is dedicated to the surface reconstruction and, instead of using the radiance fields, represents the scene implicitly by modelling the signed distance function (SDF) [50, 62, 64, 69, 73]. Unlike those approaches, we focus only on the radiance field representation and consider these methods orthogonal to ours.

Although there are some methods that train radiance fields while fine-tuning the poses [30, 57] or without known cameras [6, 65], most methods need camera poses for the reconstruction. SfM, *e.g.*, COLMAP [52, 53], is typically used for estimating the poses, which also produces a (sparse) point cloud. Our approach makes use of this by-product of the pose recovery process instead of only using the poses themselves.

Field representations. When a single MLP is used to represent the entire scene, everything is baked into a single set of parameters which cannot be modified, because any change to the parameters has a non-local effect, *i.e.*, it changes the scene at multiple unrelated places. To overcome this problem, others have experimented with different representations of the radiance fields [9, 10, 17, 31, 40, 54, 67]. A common practice is to represent the scene as a shallow MLP and an efficiently encoded voxel grid [9, 10, 31, 40], although, the encoded voxel grid can be used to represent the radiance field directly [17, 54]. The voxel grid can be encoded as a sparse tensor [17, 31, 54], a factorisation of the 4D tensor [9, 10], or a hashmap [40]. When these structures are queried, trilinear interpolation is used to combine the

¹<https://github.com/jkulhanek/tetra-nerf>

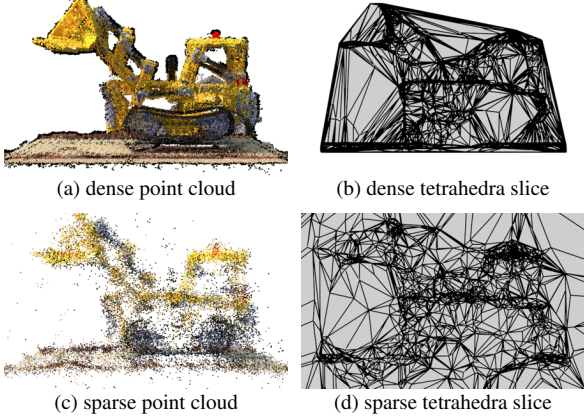


Figure 2. Input point cloud and a slice through the triangulated tetrahedra. Note that smaller tetrahedra are created closer to the surface of the scene, *i.e.*, regions close to the surface are represented with a finer resolution.

feature vectors stored in the containing voxel corners. Unfortunately, the hashmaps [40] and hierarchical representations [10] have the same non-local effect problem, and the rest of the approaches rely on subsequent upsampling of the field and can be overly complicated. Also, feature vectors cannot be placed arbitrarily in the 3D space as they must lie on the grid with a fixed resolution. In contrast, our approach is much more flexible as it stores the feature vectors freely in 3D space.

Finally, Point-NeRF [67] represents the scene as a point cloud of features which are queried using k -nearest neighbours search. However, when the point cloud contains sparse regions, the rays do not intersect any neighbourhoods of any points and the pixels stay empty without the ability to optimise. Therefore, Point-NeRF [67] relies on gradually adding more points during training and increasing the scene complexity. Since we use the triangulation of the point cloud rather than the point cloud itself, our representation is dense and does not suffer from empty regions. Therefore, we do not have to add any points during training.

3. Method

A common strategy in the literature is to represent the scene explicitly through a voxel volume. In contrast to this uniform subdivision, we investigate using an adaptive subdivision of the scene. In many scenarios, an approximation of the scene geometry is either given, *e.g.*, when using SfM to compute the input camera poses, or can be computed, *e.g.*, via MVS or single-view depth predictions [76]. This allows us to compute an adaptive subdivision of the scene via Delaunay triangulation [14] of such a point cloud. This results in a set of non-overlapping tetrahedra, where smaller tetrahedra are created close to the surface of the scene (*c.f.* Fig. 2). In the following, we explain how this adaptive subdivision of the scene can be used instead of voxels for volume rendering and neural rendering.

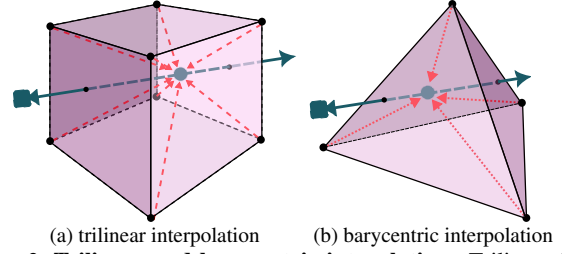


Figure 3. **Trilinear and barycentric interpolation.** Trilinear interpolation is a weighted combination of the eight voxel corners. Barycentric interpolation weights the four vertices of the tetrahedron vertices based on the barycentric coordinates [16].

3.1. Preliminaries

Neural Radiance Fields (NeRFs) [38] represent a scene through an implicit function $F(\mathbf{x}, \mathbf{d})$, often modelled via a neural network, that returns a colour value and a volume density prediction for a given 3D point \mathbf{x} in the scene observed from a viewing direction \mathbf{d} . Volume rendering [36] is used to synthesise novel views: For each pixel in a virtual view, we project a ray from the camera plane into the scene and sample the radiance field to obtain the colour and density values \mathbf{c}_i and σ_i at distances t_i along the ray, where $i = 1 \dots N$. The individual samples are then combined to predict the colour \mathbf{C} for the pixel:

$$\mathbf{C} = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \text{ where} \quad (1)$$

$$T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right),$$

and $\delta_i = t_i - t_{i-1}$ is the distance between adjacent samples. This process is fully differentiable. Thus, after computing the mean squared error (MSE) between the predicted and ground truth colours associated with each ray, we can back-propagate the gradients to the radiance field.

Voxel-based feature fields such as NSVF [31] represent the scene as a voxel grid and an MLP. Each grid point of the voxel grid is assigned a trainable vector. There are eight feature vectors associated with a single voxel, but the vectors are shared between neighbouring voxels. For each query point sampled along the ray, its corresponding voxel is found first. A feature for the point is then computed via trilinear interpolation of the features of the voxel grid points based on the position of the query point (*c.f.* Fig. 3, left). The resulting feature vector is passed through an MLP to predict the density and appearance vector. The appearance vector is combined with the ray direction and passed through a second MLP in order to compute a view-dependent colour.

3.2. Tetrahedra fields

Given a set of points in 3D space, we build the tetrahedra structure by triangulating the points. We apply the Delaunay Triangulation [14] to obtain a set of non-overlapping tetrahedra whose union is the convex hull of the original points. Fig. 2 shows example tetrahedra obtained by triangulating dense and sparse COLMAP reconstructed point clouds [52, 53]. Note that the resulting representation is adaptive as it uses a higher resolution (smaller tetrahedra) closer to the surface and larger tetrahedra to model regions farther away from the surface.

We associate all vertices of the tetrahedra with trainable vectors. As in the voxel grid case, vertices, and thus vectors, are shared between adjacent tetrahedra. The resulting tetrahedra field can be queried in the same way a voxel grid representation is queried: for each query point, we first find the tetrahedron containing the point. Instead of the trilinear interpolation used for voxel volumes, we use barycentric interpolation [16] to compute a feature vector for the query point from the four feature vectors stored at the tetrahedron’s vertices (*c.f.* Fig. 3). To this end, we compute the query point \mathbf{x} ’s barycentric coordinates λ , which express the point’s 3D coordinates as a unique weighted combination of the 3D coordinates of the tetrahedron’s vertices. In particular, the weight for a vertex is the volume of the tetrahedron constructed from the query point and the face opposite to the vertex divided by the volume of the full tetrahedron:

$$\lambda = \left(\frac{V_{x234}}{V_{1234}}, \frac{V_{1x34}}{V_{1234}}, \frac{V_{12x4}}{V_{1234}}, \frac{V_{123x}}{V_{1234}} \right), \quad (2)$$

where V_{1234} is the volume of the full tetrahedron and V_{x234} , V_{1x34} , \dots , are volumes of tetrahedra with 1st, 2nd, \dots , vertex replaced by \mathbf{x} . The same weights λ can be applied to the feature vectors of the vertices to obtain the query feature.

The interpolated features are used as the input to a small MLP in order to predict density and colour at the query point. We first pass the barycentric-interpolated features through a three-layer MLP to compute the density and appearance features. We then concatenate the appearance features with the ray direction vector encoded using Fourier features [38, 56] and pass the result through a single linear layer to get the raw RGB colour value.

Finally, to map the raw density values $\bar{\sigma}_i$ returned by the network to the volume density σ_i required by volume rendering, we apply the softplus activation function [3]. For the RGB colour value, we use the sigmoid function [17].

3.3. Efficiently querying a tetrahedra field

Determining the corresponding voxel for a given query point can be done highly efficiently via hashing [41]. In contrast, finding the corresponding tetrahedron for a query point is more complex. This in turn can significantly impact rendering, and thus training, efficiency [17, 40, 54].

In order to efficiently look up the corresponding tetrahedra, we exploit that we are not considering isolated points, but points sampled from a ray. We compute the tetrahedra that are intersected by the ray, allowing us to march through the tetrahedra rather than computing them individually per point. The relevant tetrahedra can be found efficiently using acceleration structures for fast ray-triangle intersection computations, *e.g.*, via NVidia’s OptiX library [42]: We first compute the intersections between the rays from a synthetic view and all faces of the tetrahedra. For each ray, we take the first 512 intersected triangles², and determine the corresponding tetrahedra. The tetrahedra can be ordered based on the intersections along the ray, allowing us to easily march through the tetrahedra to determine which tetrahedron to use for a given query point.

A side benefit of computing ray-triangle intersections is that we can simplify computing the barycentric coordinates: for each intersection, we compute the 2D barycentric coordinates w.r.t. the triangle. We then obtain the 3D barycentric coordinates w.r.t. the associate tetrahedron by simply adding zero for the vertex opposite to the triangle. For a query point inside a tetrahedron, we can compute its tetrahedron barycentric coordinates by linearly interpolating between the barycentric coordinates of the two intersections of the ray and the tetrahedron.

3.4. Coarse and fine sampling

We follow the common practice of having a two-stage sampling procedure [9, 38]. In the coarse stage, we sample uniformly along the ray. In the fine stage, we use the density weights from the coarse sampling stage to bias the sampling towards sampling closer to the potential surface. Following [38], we use the stratified uniform sampling for the coarse stage. The stratified uniform sampling splits the ray into equally long intervals and samples uniformly in each interval. Unlike NeRF [38], we limit the sampling to the space occupied by tetrahedra. In the fine sampling stage, we use the same network as in the coarse sampling stage.

For the fine sampling stage, we take the accumulated weights w_i from the coarse sampling:

$$\bar{w}_i = (1 - \exp(-\sigma_i \delta_i)) \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right). \quad (3)$$

These weights are the coefficients used in Equation 1 as multipliers for the colours [38]. We obtain w_i by normalizing \bar{w}_i . Following [38], we sample set of fine samples using weight w_i . To render the final colour, we merge the dense and fine samples and use all in the rendering equation [9].

²For efficiency, we only consider a fixed number of triangles. As discussed later on, this can degrade results in larger scenes / scenes with a fine-grained tetrahedralisation, where more intersections are needed. Naturally, more triangles can be considered at the cost of longer run-times.

4. Experiments

We compare Tetra-NeRF to relevant methods on the commonly used Blender dataset [38], the real-world Tanks and Temples dataset [25], and the challenging object-centric mip-NeRF 360 dataset [4]. To demonstrate its effectiveness, we compare Tetra-NeRF to a dense-grid representation and evaluate it with reduced quality of the input point cloud. We start by describing the exact hyperparameters used.

4.1. Implementation details

Generating point cloud & triangulation. Given a set of posed images, we use the COLMAP reconstruction pipeline [52, 53] to get the point cloud used in our tetrahedra field representation. We then reduce the size of the resulting point cloud such that if the number of points is larger than 10^6 , we subsample 10^6 points randomly. For the Blender dataset experiments where the number of points is smaller, we add more randomly generated points. In that case, the number of random points is half the number of original points. The reason for adding the points is that with a low number of points, some pixels on edges may not intersect any tetrahedra, potentially producing artefacts on the edges. Each added point is sampled as follows: We sample a random point x_0 from the original point cloud, sample a random normal vector n , and a number $\alpha \sim \mathcal{N}(\bar{d}, \bar{d}^2)$, where \bar{d} is the average spacing of the original point cloud, *i.e.* the average distance between each point and its six closest neighbours. We then add the point $x = x_0 + \alpha n$.

Initialization. Given the processed point cloud, we apply the Delaunay triangulation [14] to get a set of tetrahedra. To this end, we use the CGAL library [2], which in our experiments runs in order of milliseconds. Following [40], we initialise the features of size 64 at the vertices of the tetrahedra with small values sampled uniformly in the range -10^{-4} to 10^{-4} . However, to allow the model to reuse the information contained in the point cloud, we set the first four dimensions of the feature field to the RGBA colour values (rescaled to interval $[0, 1]$) stored at the associated points in the point cloud. The alpha value of all original points is one, whereas all randomly sampled points have an alpha value of zero. For the MLP, we follow the common practice of using the Kaiming uniform initialisation [19]. The hidden sizes in all MLPs are 128.

Training. During training, we sample batches of 4,096 random rays from random training images. We use volumetric rendering to predict the colour of each ray. The gradients are computed by backpropagating the MSE loss between the predicted colour and the ground truth colour value. We use the RAdam optimizer [32] and decay the learning rate exponentially from 10^{-3} to 10^{-4} in 300k steps. The training code is built on top of the Nerfstudio framework [57], and the tetrahedra field is implemented in CUDA and uses the OptiX library [42]. We train on a single NVIDIA A100

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF [38]	31.00	0.947	0.081
NSVF [31]	31.77	0.953	-
mip-NeRF [3]	34.51	0.961	0.043
instant-NGP [40]	33.18	-	-
Plenoxels [17]	31.71	0.958	0.049
Point-NeRF ^{col} [67]	31.77	0.973	0.062
Point-NeRF ^{mvs} [67]	33.31	0.978	0.049
Tetra-NeRF	32.53	0.982	0.041

Table 1. **Results on the Blender dataset** [38] averaged over all scenes in the dataset. Even though we use the same input point cloud as Point-NeRF^{col}, we outperform it greatly. We perform on par with Point-NeRF^{mvs} even though it uses many more points and densifies the point cloud during training. We highlight the **best**, **second**, and **third** values.

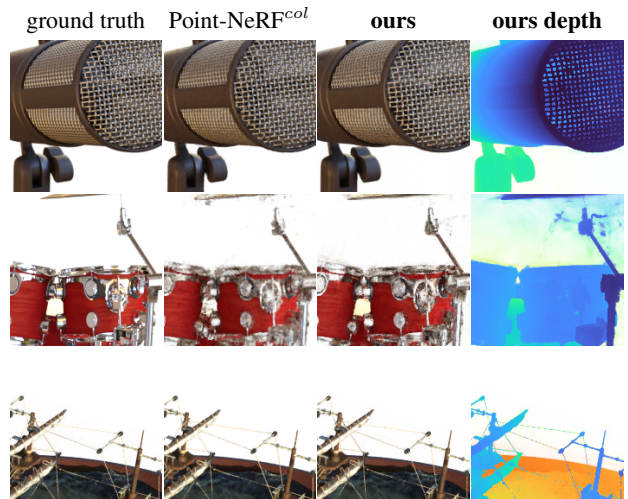


Figure 4. **Qualitative results on the Blender dataset.** We compare with Point-NeRF^{col} as we use the same input point cloud. In the **top** picture we can see that Tetra-NeRF is able to represent fine details well on the *mic* scene. On the *drums* scene (**middle**), both methods struggle with shiny materials, but our method performs slightly better. On the **bottom**, Tetra-NeRF can render thin ropes.

GPU and the training speed ranges from 15k rays per second to 30k rays per second. The speed depends on how well-structured the triangulation is, how many vertices there are, and if there is empty space around the object. The full training with 300k iterations takes between 11 and 24 hours, depending on the scene complexity. However, good results are typically obtained much earlier, *e.g.*, in 100k iterations.

4.2. Evaluation

Blender dataset [38] results. We compare to relevant baselines on the standard Blender dataset. We used the same split and evaluation procedure as in the original NeRF paper [38] and the same SSIM implementation as Point-NeRF [67]. In order to ensure a fair comparison with Point-NeRF [67] when COLMAP points were used, we use the exact same COLMAP reconstruction as Point-NeRF. We re-

	<i>hotdog</i>		<i>ship</i>	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Point-NeRF ^{static}	29.91	0.978	19.35	0.905
Tetra-NeRF	33.31	0.989	31.13	0.994

Table 2. **Comparison with Point-NeRF with disabled point cloud growing and pruning** shows that Point-NeRF performs significantly worse in all measured metrics because it cannot handle the sparse point cloud.



Figure 5. **Comparison with dense grid field representation.** Tetra-NeRF on the **left**, dense grid on the **right**. We used the same number of parameters for both methods, however, thanks to the adaptive nature of tetrahedra fields, Tetra-NeRF produces significantly better rendering as it is able to focus on relevant parts of the scene.

port the PSNR, SSIM, and LPIPS (Vgg) [75] metrics. Tab. 1 shows averaged results, Fig. 4 shows qualitative results, and the *Supp. Mat.* shows results for individual scenes.

When we use the exact same COLMAP points as Point-NeRF (row *Point-NeRF^{col}*), we outperform it in all three metrics. We score comparably with *Point-NeRF^{mv}*, even though it uses a much higher quality of the initial point cloud, which it generates from a jointly trained model. Also note, that both of these Point-NeRF configurations grow the point cloud during training and, therefore, the complexity of the scene representation grows. For us, the points are fixed and the number of parameters stays the same. We also outperform Plenoxels [17], which uses a sparse grid. Note, that same as Point-NeRF, Plenoxels also gradually increases the representation complexity by subdividing the grid resolution at predefined training epochs. Even though both mip-NeRF and instant-NGP outperform our approach in terms of PSNR, our method is slightly better in terms of SSIM and on par with mip-NeRF in terms of LPIPS.

To analyze the tetrahedra field representation, we compare it with the point cloud field representation used in Point-NeRF. In Table 2, we compare our method to Point-NeRF when we disable point cloud growth and pruning. We show the results on two scenes from the Blender dataset [38] which were selected in the Point-NeRF paper. With this setup, we vastly outperform Point-NeRF in all metrics. The reason is that Point-NeRF requires the point cloud to be dense such that all rays have a chance to intersect a neighbourhood of a point. Since we use a continuous representation (tetrahedra field) rather than a discrete one, we can achieve good results even for sparser point clouds.

Comparison with the dense grid representation. In order

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NV [33]	23.70	0.848	0.260
NeRF [38]	25.78	0.864	0.198
NSVF [31]	28.40	0.900	0.153
Point-NeRF [67]*	28.35	0.942	0.090
Tetra-NeRF	28.90	0.957	0.059

Table 3. **Result on the Tanks and Temples dataset [25]** as processed by NSVF [31] averaged over all scenes. We show the PSNR, SSIM, and LPIPS (Alex) [75] metrics. We highlight the **best**, **second**, and **third** values. *Note, that Point-NeRF [67] results differ from the paper as they were recomputed with the resolution used in other methods.

to show the utility of the adaptive tetrahedra field representation, we compare it to a dense grid representation. Similarly to NSVF [31], we split the 3D scene bounding box into equally-sized voxels. When querying the field, we find the voxel to which the query point belongs and perform trilinear interpolation of the eight corners of the voxel. To ensure a fair comparison, we choose the grid resolution such that the number of grid points is the lowest cube number larger than the number of points of the original point cloud. All other hyperparameters are kept the same as for Tetra-NeRF. We evaluate both approaches on two scenes from the Blender dataset [38]. The dense grid representation only scores PSNRs 18.81 and 18.91 on the *lego* and *ship* scenes respectively, whereas Tetra-NeRF scores PSNRs 33.79 and 30.69 respectively. The results can be seen in Figure 5. Note, that in this experiment we only trained for 100,000 iterations to save computation time. From the numbers and the figure, we can clearly see that the dense grid resolution is not sufficient to reconstruct the scene in enough detail. Because Tetra-NeRF uses an adaptive subdivision, which is more detailed around the surface, we are better able to focus on relevant scene parts. With the same number of trainable parameters, we thus obtain better results.

Tanks and Temples [25] dataset. To be able to compare with Point-NeRF [67] on real-world data, we have evaluated Tetra-NeRF on the Tanks and Temples [25] dataset. We use the same setup as in NSVF [31], where the object is masked. We report the usual metrics: PSNR, SSIM, LPIPS (Alex)³ We used the dense COLMAP reconstruction to get the point cloud used in the tetrahedra field. The quantitative results are shown in Table 3, the qualitative results are presented in Figure 6, and the per-scene results are given in *Supp. Mat.* Note, that the results originally reported in the Point-NeRF paper [67] were evaluated with a resolution different from the compared methods. Therefore, to ensure a fair comparison, we have recomputed the metrics in the publicly available Point-NeRF’s predictions with the full resolution of 1920×1080 as used in NSVF [31]. In the public dataset, one of the scenes had corrupted camera pa-

³We always choose the type of LPIPS (Alex or Vgg) such that we can compare with more methods as some only evaluate using one type.

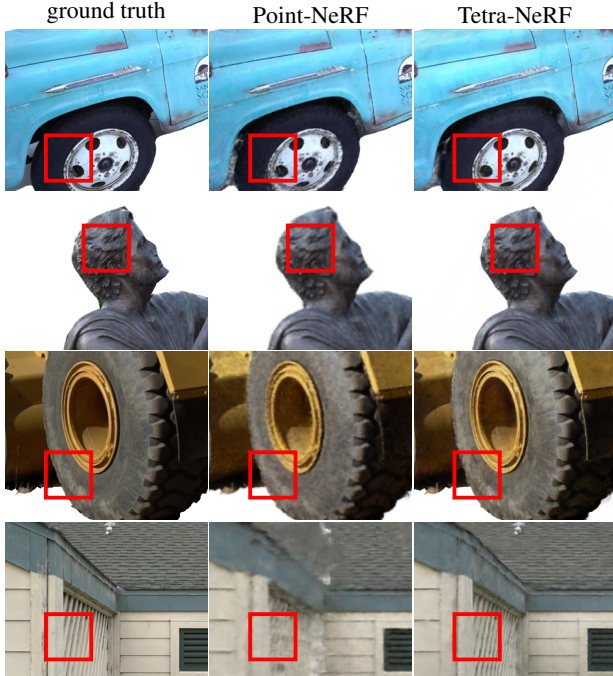


Figure 6. **Results on Tanks and Temples dataset.** In the **top row**, we can see that we are able to represent the rim of the wheel better than Point-NeRF. Similarly, we can represent the tyre better in the **third row**. Finally, in the **bottom**, Point-NeRF fails on the wall and the roof, whereas Tetra-NeRF can render these parts well.

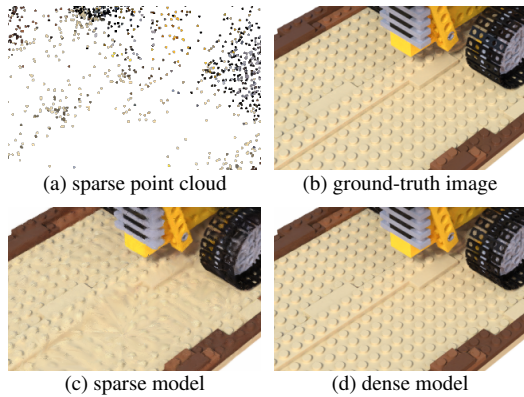


Figure 7. **Sparse vs. dense model comparison detail.** Even though the area in the bottom contains few points, we are still able to reconstruct at least low-frequency data.

rameters and we had to reconstruct the poses again for the Ignatius scene. With the reconstructed poses, our method performs slightly worse on that scene compared to others.

The results show that our method outperforms the baselines in all compared metrics. This indicates that even though Point-NeRF relies on the ability to grow the point cloud density during training, this is not needed when using a continuous instead of a discrete representation.

Sparse and dense point cloud comparison. Previous experiments have used dense COLMAP reconstructions. However, dense point clouds may not be required to achieve

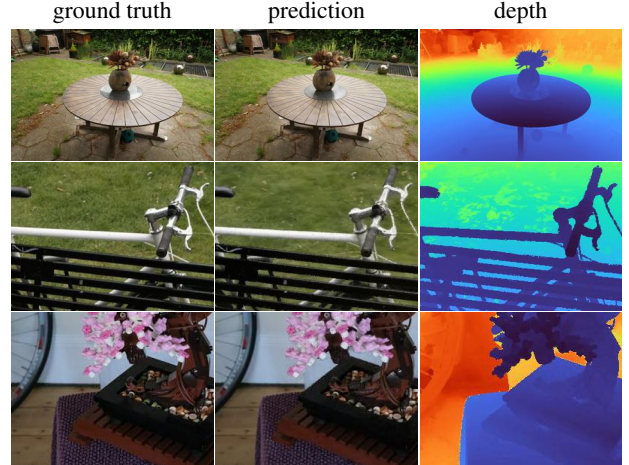


Figure 8. **Mip-NeRF 360 results.** **Top row:** we show the results on the garden scene. **Middle row:** we can see that Tetra-NeRF have some problems rendering grass. **Bottom row:** Tetra-NeRF is able to represent the delicate texture of the tablecloth well.

good reconstructions. In this section, we compare models trained on dense reconstructions and sparse reconstructions. Table 5 compares PSNR, SSIM metrics on two synthetic scenes from the Blender dataset [38] and two real-world scenes from the mip-NeRF 360 dataset [4]. We also show the number of vertices of the tetrahedra which is the size of the input point cloud, including the random points added at the beginning. Note, that in this experiment we only trained for 100,000 iterations to save computation time.

From the results, we can see, that for the Blender dataset [38] case, the dense model is much better. This is to be expected since the number of points obtained from the sparse reconstruction on the Blender dataset is very small and our MLP does not have enough capacity to represent fine details. However, even in regions with zero point coverage, we are still able to provide at least some low-frequency data (*c.f.* Fig. 7). On the real-world scenes, the sparse point cloud model is almost on par with the dense one. The sparse reconstructions on the real-world dataset provide many more points compared to the synthetic dataset case and the coverage is sufficient. From these results, we conclude that on real-world data the sparse model can be sufficient to achieve a good performance.

Mip-NeRF 360 [4] dataset. We have further evaluated our method on the mip-NeRF 360 [4] dataset. In order to ensure a fair comparison with mip-NeRF 360 [4], we trained and evaluated on four times downsized images for the outdoor scenes, and two times downsized images for the indoor scenes. The quantitative results are presented in Table 4 and the qualitative results are shown in Figure 8.

From the results we can see, that Tetra-NeRF is able to outperform both the vanilla NeRF [38] and mip-NeRF [3]. We also outperform Stable View Synthesis [49] in terms of PSNR and score comparably in terms of LPIPS. This is pos-

	<i>Outdoor</i>			<i>Indoor</i>			<i>Mean</i>		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF [15, 38]	21.46	0.458	0.515	26.84	0.790	0.370	23.85	0.605	0.451
mip-NeRF [3]	21.69	0.471	0.505	26.98	0.798	0.361	24.04	0.616	0.441
NeRF++ [74]	22.76	0.548	0.427	28.05	0.836	0.309	25.11	0.676	0.375
Deep Blending [20]	21.54	0.524	0.364	26.39	0.844	0.261	23.70	0.666	0.318
Point-Based Neural Rendering [26]	21.66	0.612	0.302	26.28	0.887	0.191	23.71	0.734	0.253
Stable View Synthesis [49]	23.01	0.662	0.253	28.22	0.907	0.160	25.33	0.771	0.211
mip-NeRF 360 [4]	23.72	0.687	0.282	29.43	0.911	0.181	26.26	0.786	0.237
Tetra-NeRF	23.17	0.586	0.298	30.21	0.881	0.103	26.30	0.717	0.211

Table 4. **Mip-NeRF 360 dataset results.** We show the PSNR, SSIM, and LPIPS (Alex) [75] on two categories of mip-NeRF 360 [4] scenes: *outdoor*, and *indoor*. On the *outdoor* scenes, we outperform the Stable View Synthesis [49] and are comparable to mip-NeRF 360 [4], even though our method does not implement the improvements suggested in NeRF++ [74] and mip-NeRF 360 [4], and is more comparable to vanilla NeRF [38]. We highlight the **best**, **second**, and **third** values.

	<i>sparse</i>			<i>dense</i>		
	PSNR/SSIM	#points		PSNR/SSIM	#points	
Blender/lego	29.77/0.959	25,784		33.79/0.985	302,781	
Blender/ship	26.90/0.899	7,152		30.69/0.942	321,861	
360/bonsai	28.12/0.902	413,226		28.34/0.902	1,000,000	
360/garden	24.79/0.806	227,532		25.41/0.838	1,000,000	

Table 5. **Dense and sparse point cloud comparison.** We present results on two scenes from Blender [38] and mip-NeRF 360 [4] dataset. We also show the number of vertices of tetrahedra field – including the randomly sampled points.



(a) **Failure case 1:** low pointcloud density on the grass (b) **Failure case 2:** too many intersected tetrahedra

Figure 9. **Failure cases.** Left: barn scene from [25]. Right: tree hill scene from [4].

sible because outdoor scenes contain a lot of high-resolution geometries such as leaves and Stable View Synthesis is not able to overcome noise in the approximate geometry. Tetra-NeRF does not suffer from these problems as we aggregate features in 3D along rays rather than on the surface. Mip-NeRF 360 scores comparably to Tetra-NeRF in terms of PSNR and LPIPS, but it outperforms Tetra-NeRF in terms of SSIM. However, mip-NeRF 360 implements some tricks designed to boost its performance. On the other hand, our approach is based on vanilla NeRF, and we believe that the performance can be boosted similarly.

5. Limitations

One drawback of our approach is that the quality of different regions of the rendered scene depends on the density of the point cloud in these regions. If the original point cloud was constructed using 2D feature matching, there can be regions with a very low number of 3D points. We show such an example in Fig. 9(a). There are few 3D points on

the ground, resulting in blurry renderings for that region.

Another issue is that the current implementation has a limit on the number of intersected tetrahedra per ray. For large or badly structured scenes, where rays intersect many tetrahedra, this can limit the reconstruction quality. Fig. 9(b) shows such a case on a mip-NeRF 360 scene [4]. This problem can be addressed by increasing the limit on the number of visited tetrahedra, at the cost of higher memory requirements and run-time. Alternatively, coarse-to-fine schemes that start with a coarser tetrahedralisation and prune the space could potentially handle this issue.

6. Conclusion

This paper proposes a novel radiance field representation that, compared to standard voxel-based representations, is easily able to adapt to 3D geometry priors given in the form of a (sparse) point cloud. Our approach elegantly combines concepts from geometry processing (Delaunay triangulation) and triangle-based rendering (ray-triangle intersections) with modern neural rendering approaches. The representation has a naturally higher resolution in the space near surfaces, and the input point cloud provides a straightforward way to initialise the radiance field. Compared to Point-NeRF, a state-of-the-art point cloud-based radiance field representation which uses the same input as our approach, Tetra-NeRF shows clearly better results. Our method performs comparably to state-of-the-art MLP-based methods. The results demonstrate that Tetra-NeRF is an interesting alternative to existing radiance field representations that is worth further investigation. Interesting research directions include adaptive refinement and pruning of the tetrahedralisation and exploiting the fact that the surface of the scene is likely close to some of the triangles in the scene.

Acknowledgements This work was supported by the European Regional Development Fund under the project IMPACT (reg. no. CZ.02.1.01/0.0/0.0/15 003/0000468), the Czech Science Foundation (GA ČR) EXPRO (grant no. 23-07973X), and the Grant Agency of the Czech Technical University in Prague (grant no. SGS22/112/OHK3/2T/13).

Supplementary Material

First, in Section A, we describe the **attached video** where we illustrate the tetrahedra field and the optimisation process. Next, we extend Sections 4.2, 4.3, and 4.4 from the main paper by giving detailed results on the Blender [38], Tanks and Temples [25], and mip-NeRF 360 [4] datasets in Sections B, C, and D. Finally, we evaluate how the performance changes when varying the number of input points in Section E.

A. Attached video

To present the idea of the paper visually, we include a video illustrating the tetrahedra field representation and showing the results of the optimisation at different points in the early stages of the training. In the video, we have used the *garden* scene from the mip-NeRF 360 dataset [4]. We show the optimisations using both the sparse and dense input point clouds. The video starts by showing the initial point cloud and the triangulated tetrahedra. It then shows how the scene is optimised from the first iteration. Finally, it presents the resulting video generated from the fully trained model.

	PSNR↑								
	<i>chair</i>	<i>drums</i>	<i>ficus</i>	<i>hotdog</i>	<i>lego</i>	<i>materials</i>	<i>mic</i>	<i>ship</i>	<i>mean</i>
NPBG [1]	26.47	21.53	24.60	29.01	24.84	21.58	26.62	21.83	24.56
NeRF [38]	33.00	25.01	30.13	36.18	32.54	29.62	32.91	28.65	31.01
NSVF [31]	33.19	25.18	31.23	37.14	32.54	32.68	34.27	27.93	31.77
Mip-NeRF [3]	37.14	27.02	33.19	39.31	35.74	32.56	38.04	33.08	34.51
Instant-NGP [40]	35.00	26.02	33.51	37.40	36.39	29.78	36.22	31.10	33.18
Plenoxels [17]	33.98	25.35	31.83	36.43	34.10	29.14	33.26	29.62	31.71
Point-NeRF ^{col} [67]	35.09	25.01	33.24	35.49	32.65	26.97	35.54	30.18	31.77
Point-NeRF ^{mvs} [67]	35.40	26.06	36.13	37.30	35.04	29.61	35.95	30.97	33.31
Tetra-NeRF	35.05	25.01	33.31	36.16	34.75	29.30	35.49	31.13	32.53

	SSIM↑								
	<i>chair</i>	<i>drums</i>	<i>ficus</i>	<i>hotdog</i>	<i>lego</i>	<i>materials</i>	<i>mic</i>	<i>ship</i>	<i>mean</i>
NPBG [1]	0.939	0.904	0.940	0.964	0.923	0.887	0.959	0.866	0.923
NeRF [38]	0.967	0.925	0.964	0.974	0.961	0.949	0.980	0.856	0.947
NSVF [31]	0.968	0.931	0.973	0.980	0.960	0.973	0.987	0.854	0.953
Mip-NeRF [3]	0.981	0.932	0.980	0.982	0.978	0.959	0.991	0.882	0.961
Plenoxels [17]	0.977	0.933	0.890	0.985	0.976	0.975	0.980	0.949	0.958
Point-NeRF ^{col} [67]	0.990	0.944	0.989	0.986	0.983	0.955	0.993	0.941	0.973
Point-NeRF ^{mvs} [67]	0.991	0.954	0.993	0.991	0.988	0.971	0.994	0.942	0.978
Tetra-NeRF	0.990	0.947	0.989	0.989	0.987	0.968	0.993	0.994	0.982

	LPIPS↓								
	<i>chair</i>	<i>drums</i>	<i>ficus</i>	<i>hotdog</i>	<i>lego</i>	<i>materials</i>	<i>mic</i>	<i>ship</i>	<i>mean</i>
NPBG [1]	0.085	0.112	0.078	0.075	0.119	0.134	0.060	0.210	0.109
NeRF [38]	0.046	0.091	0.044	0.121	0.050	0.063	0.028	0.206	0.081
Mip-NeRF [3]	0.021	0.065	0.020	0.027	0.021	0.040	0.009	0.138	0.043
Plenoxels [17]	0.031	0.067	0.026	0.037	0.028	0.057	0.015	0.134	0.049
Point-NeRF ^{col} [67]	0.026	0.099	0.028	0.061	0.031	0.100	0.019	0.134	0.062
Point-NeRF ^{mvs} [67]	0.023	0.078	0.022	0.037	0.024	0.072	0.014	0.124	0.049
Tetra-NeRF	0.016	0.073	0.023	0.027	0.022	0.056	0.011	0.103	0.041

Table 6. **Detailed results on the Blender dataset [38].** We show the PSNR, SSIM, and LPIPS (Vgg) results averaged over the testing images. We highlight the **best**, **second**, and **third** values. We outperform Point-NeRF^{col} [67], which was evaluated with the same input point cloud as our method.

B. Blender results

In Sec. 4.2 and Tab. 1 of the main paper, we presented results averaged over all scenes in the Blender dataset [38]. In the following, we present results individually per scene. The quantitative results can be seen in Table 6. We report the PSNR, SSIM, and LPIPS (Vgg) [75] metrics. The evaluation protocol is the same as in Point-NeRF [67], and we use the same input point cloud as Point-NeRF^{col}. From the results, we can see that we outperform the closest approach to ours, Point-NeRF, in all three metrics on almost all scenes when using the same input point cloud (row Point-NeRF^{col}). We are slightly outperformed by Point-NeRF^{mvs}, which uses denser input point clouds obtained from its end-to-end optimised Multi-View Stereo (MVS) pipeline. The results are more pronounced on the *figus* scene where Point-NeRF performs exceptionally well, outperforming other baselines, including mip-NeRF [3]. Note that both Point-NeRF configurations grow the point cloud during training and, therefore, the complexity of the scene representation grows. For us, the points are fixed, and the number of parameters stays the same. In most scenes, we also outperform Plenoxels [17], which uses a sparse grid. Note that same as Point-NeRF, Plenoxels also gradually increases the representation complexity by subdividing the grid resolution at predefined training epochs. Even though mip-NeRF [3] outperforms our approach in terms of PSNR, our method is slightly better in terms of SSIM and on par with mip-NeRF [3] in terms of LPIPS.

We also show rendered images from all Blender scenes in Figures 10 and 11. Some artefacts can only be noticed on scenes with highly reflective surfaces – *materials*, and *drums*. By closely inspecting the produced depth maps, one can notice that sometimes the density is non-zero in large tetrahedra connecting different parts of the object. A possible cause could be the combination of the training process and the implicit bias of our model. Since the tetrahedra field uses barycentric interpolation, the features will change linearly in the tetrahedra connecting different parts of the object. However, these tetrahedra should have a density of zero everywhere except for the regions close to the vertices. For the shallow MLP, it is difficult to represent such a function, and there is not enough pressure in the optimisation process to enforce it because the error will be close to zero since the background is white, without any texture, independently of how the density is distributed in these regions.



Figure 10. **Results on the Blender dataset [38] (part 1).** We show the **ground-truth** image, the **prediction**, and the **predicted depth map** on scenes: *ship* (top), *lego*, *mic*, and *chair* (bottom). The red squares highlight regions where there are visible errors in the images. Notice how our approach is able to recover fine textures in *chair* scene and tiny details in the *ship* scene geometry. However, the density seems to be non-zero in tetrahedra connecting different legs of the *chair*.

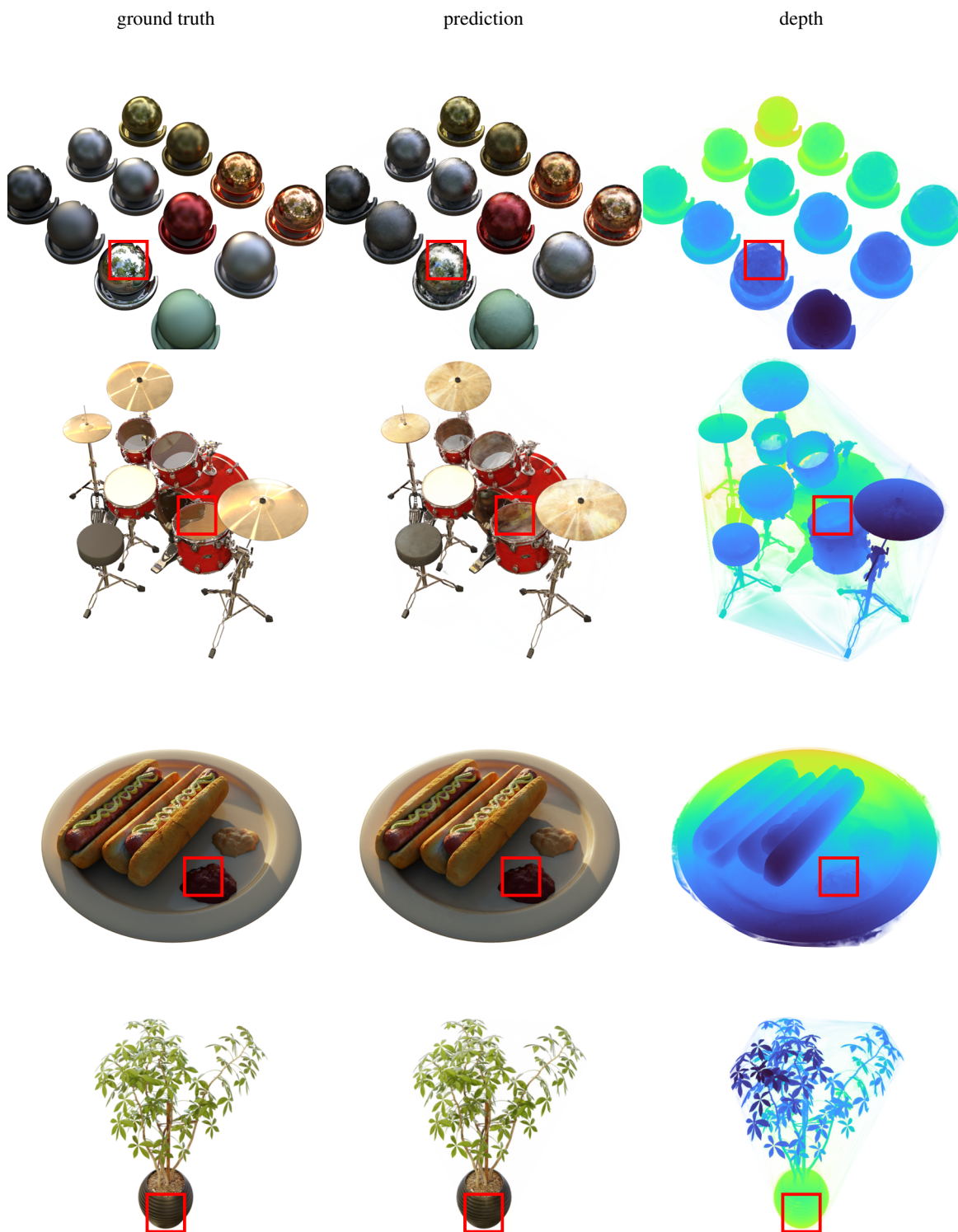


Figure 11. **Results on the Blender dataset [38] (part 2).** We show the **ground-truth** image, the **prediction**, and the **predicted depth map** on scenes: *materials* (top), *drums*, *hotdog*, and *figus* (bottom). The red squares highlight regions where there are visible errors in the images. Some artefacts can be noticed in the top two scenes, where there are highly reflective surfaces. Also, the density seems to be non-zero in tetrahedra connecting distant parts of the 3D object.

C. Tanks and Temples results

We show the detailed per-scene results for the Tanks and Temples dataset [25]. Note that to be able to compare with the Point-NeRF method [67], we used the data pre-processing and split proposed by NSVF [31], where the background is masked out. The quantitative results can be seen in Table 7. We report the PSNR, SSIM, and LPIPS (Alex) [75] metrics. The evaluation protocol is the same as in Point-NeRF [67], but we evaluate it on the original resolution, the same as other compared methods. To be able to compare with Point-NeRF [67] which evaluated on a lower-resolution images⁴, we have recomputed its metrics with the same full image resolution 1920×1080 . To obtain the initial point clouds, we used the dense COLMAP reconstruction, which was computed using the known intrinsic and extrinsic camera parameters. However, the NSVF [31] published split had corrupted camera parameters for the *ignatius* scene, and we had to run COLMAP reconstruction from scratch to obtain both the camera poses and intrinsics. This is likely the reason for the worse results on that scene. Otherwise, we outperform all baseline methods on all metrics.

We also extend Fig. 6 from the main paper and show more qualitative results in Figure 12. Compared to Point-NeRF [67], our method produces less noisy images. When looking at the depth maps, we can observe similar non-zero density artefacts to the ones observed on the Blender dataset. Same as in the Blender dataset case, a likely cause could be a combination of the implicit bias of our method and the usage of non-textured background.

PSNR \uparrow						
	<i>barn</i>	<i>caterpillar</i>	<i>family</i>	<i>ignatius</i>	<i>truck</i>	<i>mean</i>
NV [33]	20.82	20.71	28.72	26.54	21.71	23.70
NeRF [38]	24.05	23.75	30.29	25.43	25.36	25.78
NSVF [31]	27.16	26.44	33.58	27.91	26.92	28.40
Point-NeRF [67]*	27.40	25.58	33.57	28.39	26.83	28.35
Tetra-NeRF	28.86	26.64	34.27	27.17 [†]	27.58	28.90

SSIM \uparrow						
	<i>barn</i>	<i>caterpillar</i>	<i>family</i>	<i>ignatius</i>	<i>truck</i>	<i>mean</i>
NV [33]	0.721	0.819	0.916	0.992	0.793	0.848
NeRF [38]	0.750	0.860	0.932	0.920	0.860	0.864
NSVF [31]	0.823	0.900	0.954	0.930	0.895	0.900
Point-NeRF [67]*	0.908	0.927	0.976	0.959	0.939	0.942
Tetra-NeRF	0.942	0.944	0.985	0.962	0.952	0.957

LPIPS \downarrow						
	<i>barn</i>	<i>caterpillar</i>	<i>family</i>	<i>ignatius</i>	<i>truck</i>	<i>mean</i>
NV [33]	0.117	0.312	0.479	0.280	0.111	0.260
NeRF [38]	0.111	0.192	0.395	0.196	0.098	0.198
NSVF [31]	0.106	0.148	0.307	0.141	0.063	0.153
Point-NeRF [67]*	0.142	0.118	0.034	0.064	0.091	0.090
Tetra-NeRF	0.087	0.077	0.021	0.050	0.062	0.059

Table 7. **Tanks and Temples results.** We show the PSNR, SSIM, and LPIPS (Alex) results averaged over the testing images. We highlight the **best**, **second**, and **third** values. We outperform all compared methods in all metrics on all scenes except for *ignatius*, where we did not have the correct camera parameters and had to run camera pose estimation prior training[†].

⁴<https://github.com/Xharlie/pointnerf/issues/62>

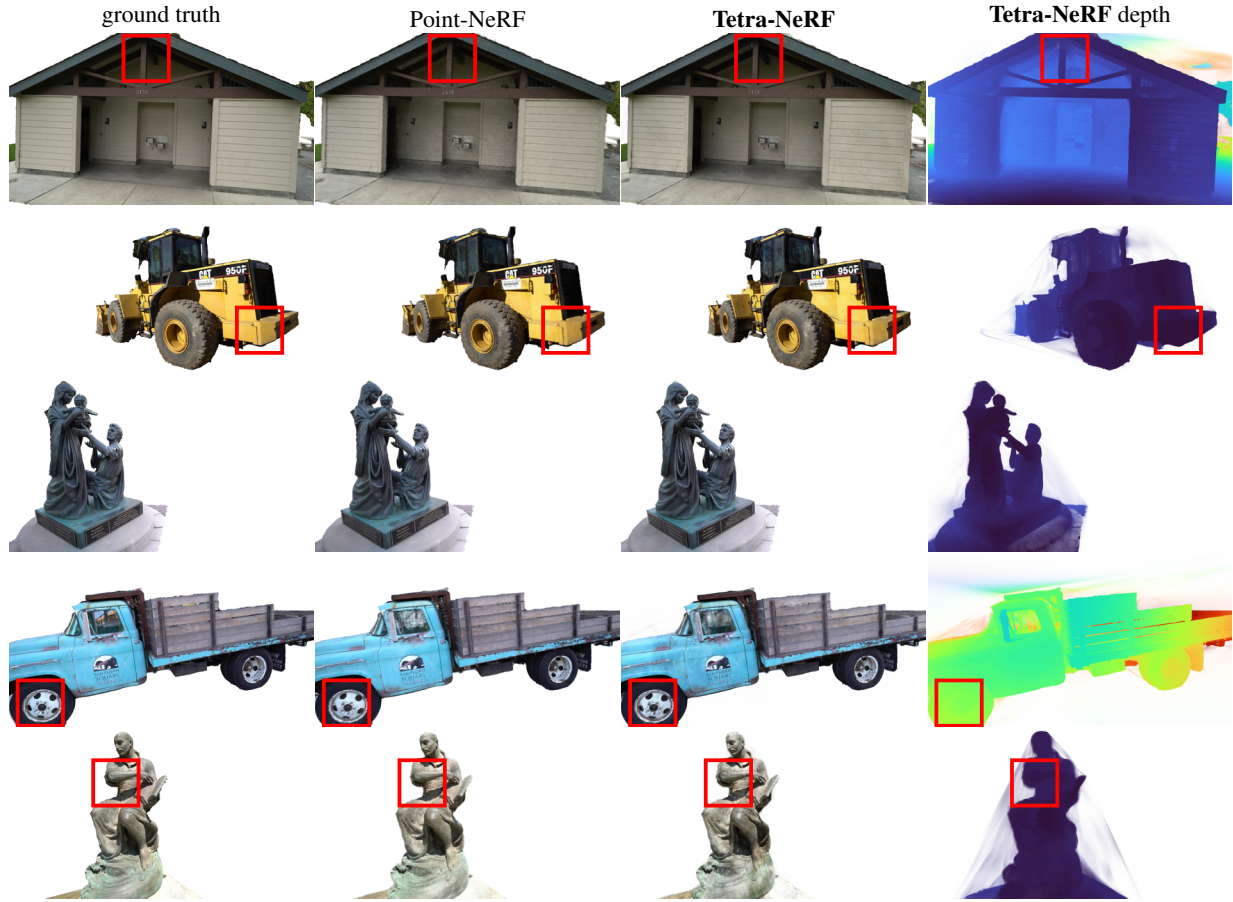


Figure 12. **Results on the Tanks and Temples dataset [25].** We show the **ground-truth** image, the **prediction**, and the **predicted depth map** on scenes: *barn* (top), *caterpillar*, *family*, *truck*, and *ignatius* (bottom). We also show a comparison with Point-NeRF [67]. The red squares highlight regions where our method has less artefacts compared to Point-NeRF. Notice how our method produces less noisy images compared to Point-NeRF. We can observe similar non-zero density artefacts to the ones observed on the Blender dataset. Again, we attribute it to a combination of implicit bias and non-textured background.

D. Mip-NeRF 360 results

We also extend the results on the mip-NeRF 360 dataset [4], presented in Table 8 and Figure 4 in the main paper, by showing detailed results for each scene. We follow the same training and evaluation procedure as mip-NeRF 360 [4], and for the outdoor and indoor scenes, we train and evaluate with 4x and 2x downsampled images, respectively. The PSNR, SSIM, and LPIPS (Alex) [75] results are presented in Table 8. In terms of SSIM, our approach performs slightly worse than Point-Based Neural Rendering [26], Stable View Synthesis [49], and mip-NeRF 360 [4]. In terms of the PSNR and LPIPS metrics, our approach performs comparable or better than these baselines. *E.g.*, the state-of-the-art mip-NeRF 360 [4] has a slightly better PSNR and Tetra-NeRF has a slightly better LPIPS. We typically outperform Stable View Synthesis [49] and competitors other than mip-NeRF 360 in terms of PSNR. Note, that similarly to Tetra-NeRF, Stable View Synthesis uses a geometric prior as input – a mesh instead of a point cloud.

We also show qualitative results for indoor and outdoor mip-NeRF 360 scenes in Figures 13 and 14. We notice more artefacts in the outdoor scenes compared to indoor ones. Small, high-frequency details such as grass are not represented well, which can be visible, *e.g.*, in the *flowers* scene. One potential cause for this behaviour, which we plan to investigate in future work, could be that the poses estimated outdoors (where the camera is typically farther away from the scene than indoors) are noisier, which makes it harder to recover fine details. For the indoor scenes, our model is able to represent the scenes with very high fidelity, including very fine details such as texts on products in the *counter* scene.

	PSNR									
	Outdoor					Indoor				
	<i>bicycle</i>	<i>flowers</i>	<i>garden</i>	<i>stump</i>	<i>treehill</i>	<i>room</i>	<i>counter</i>	<i>kitchen</i>	<i>bonsai</i>	
NeRF [15, 38]	21.76	19.40	23.11	21.73	21.28	28.56	25.67	26.31	26.81	
mip-NeRF [3]	21.69	19.31	23.16	23.10	21.21	28.73	25.59	26.47	27.13	
NeRF++ [74]	22.64	20.31	24.32	24.34	22.20	28.87	26.38	27.80	29.15	
Deep Blending [20]	21.09	18.13	23.61	24.08	20.80	27.20	26.28	25.02	27.08	
Point-Based Neural Rendering [26]	21.64	19.28	22.50	23.90	20.98	26.99	25.23	24.47	28.42	
Stable View Synthesis [49]	22.79	20.15	25.99	24.39	21.72	28.93	26.40	28.49	29.07	
mip-NeRF 360 [4]	23.95	21.60	25.09	25.98	21.99	28.24	28.40	30.81	30.27	
Tetra-NeRF	23.53	20.36	26.15	24.42	21.41	32.02	28.02	29.66	31.13	

	SSIM									
	Outdoor					Indoor				
	<i>bicycle</i>	<i>flowers</i>	<i>garden</i>	<i>stump</i>	<i>treehill</i>	<i>room</i>	<i>counter</i>	<i>kitchen</i>	<i>bonsai</i>	
NeRF [15, 38]	0.455	0.376	0.546	0.453	0.459	0.843	0.775	0.749	0.792	
mip-NeRF [3]	0.454	0.373	0.543	0.517	0.466	0.851	0.779	0.745	0.818	
NeRF++ [74]	0.526	0.453	0.635	0.594	0.530	0.852	0.802	0.816	0.876	
Deep Blending [20]	0.466	0.320	0.675	0.634	0.523	0.868	0.856	0.768	0.883	
Point-Based Neural Rendering [26]	0.608	0.487	0.735	0.651	0.579	0.887	0.868	0.876	0.919	
Stable View Synthesis [49]	0.663	0.541	0.818	0.683	0.606	0.905	0.886	0.910	0.925	
mip-NeRF 360 [4]	0.687	0.582	0.800	0.745	0.619	0.907	0.890	0.916	0.932	
Tetra-NeRF	0.614	0.470	0.775	0.613	0.456	0.894	0.850	0.877	0.905	

	LPIPS									
	Outdoor					Indoor				
	<i>bicycle</i>	<i>flowers</i>	<i>garden</i>	<i>stump</i>	<i>treehill</i>	<i>room</i>	<i>counter</i>	<i>kitchen</i>	<i>bonsai</i>	
NeRF [15, 38]	0.536	0.529	0.415	0.551	0.546	0.353	0.394	0.335	0.398	
mip-NeRF [3]	0.541	0.535	0.422	0.490	0.538	0.346	0.390	0.336	0.370	
NeRF++ [74]	0.455	0.466	0.331	0.416	0.466	0.335	0.351	0.260	0.291	
Deep Blending [20]	0.377	0.476	0.231	0.351	0.383	0.266	0.258	0.246	0.275	
Point-Based Neural Rendering [26]	0.313	0.372	0.197	0.303	0.325	0.216	0.209	0.160	0.178	
Stable View Synthesis [49]	0.243	0.317	0.137	0.281	0.286	0.182	0.168	0.125	0.164	
mip-NeRF 360 [4]	0.296	0.343	0.173	0.258	0.338	0.208	0.206	0.129	0.182	
Tetra-NeRF	0.271	0.378	0.136	0.274	0.429	0.104	0.127	0.098	0.084	

Table 8. **Detailed mip-NeRF 360 [4] results.** We show the PSNR, SSIM, and LPIPS (Alex) results averaged over the testing images. We highlight the **best**, **second**, and **third** values. Our method has a worse SSIM. However, Tetra-NeRF seems to perform comparably to mip-NeRF 360 [4] on most scenes in terms of PSNR and LPIPS, where we seem to achieve a slightly higher LPIPS and a slightly lower PSNR. We also outperform Stable View Synthesis [49] and all competitors other than mip-NeRF 360 on all scenes in terms of PSNR.

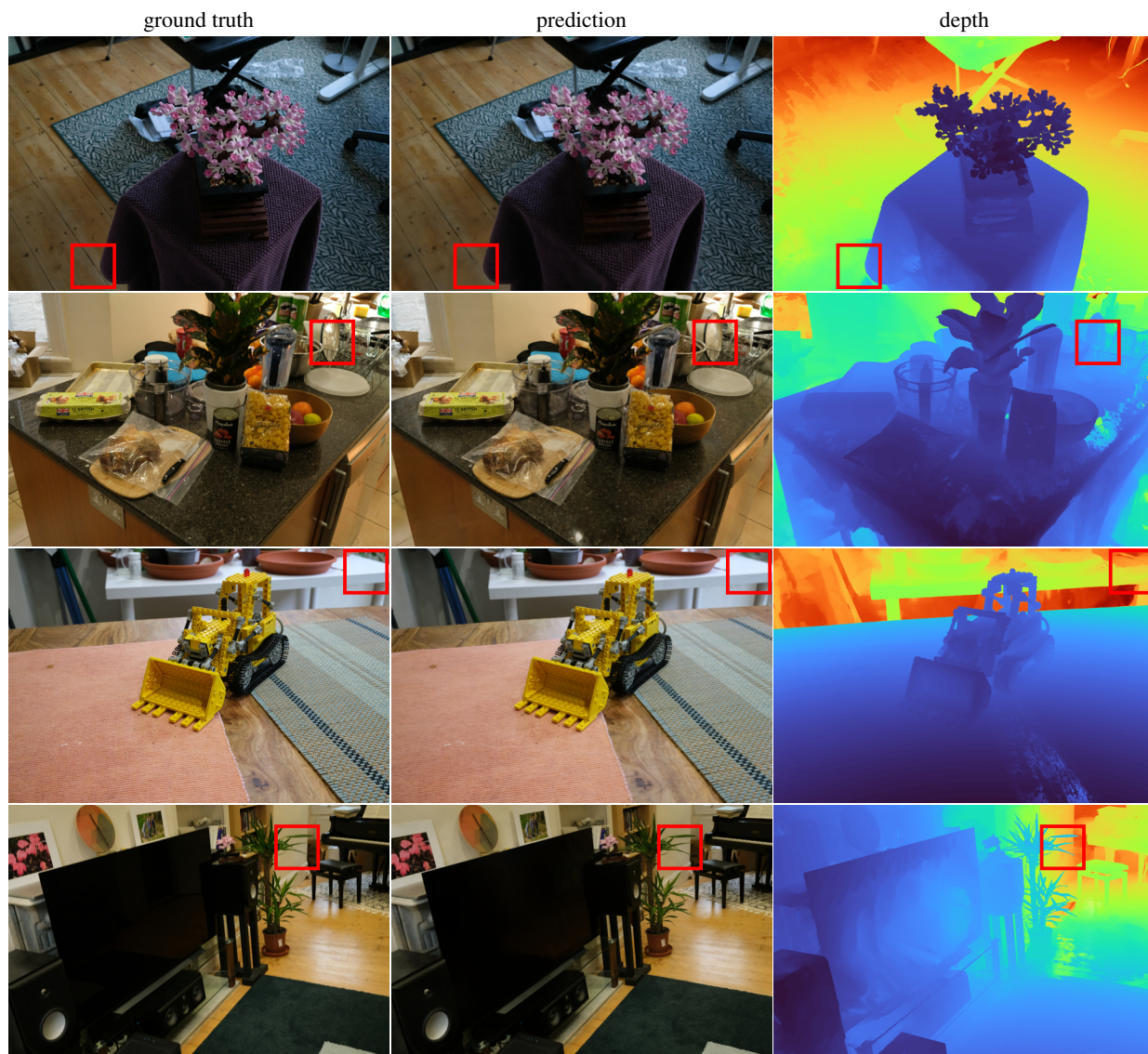


Figure 13. Results on the indoor scenes from mip-NeRF 360 dataset [4]. We show the **ground-truth** image, the **prediction**, and the **predicted depth map** on scenes: *bonsai* (top), *counter*, *kitchen*, and *room* (bottom). The red squares highlight regions where there are visible errors in the images. Tetra-NeRF is able to represent even fine details, such as texts on products in the *counter* scene, well.

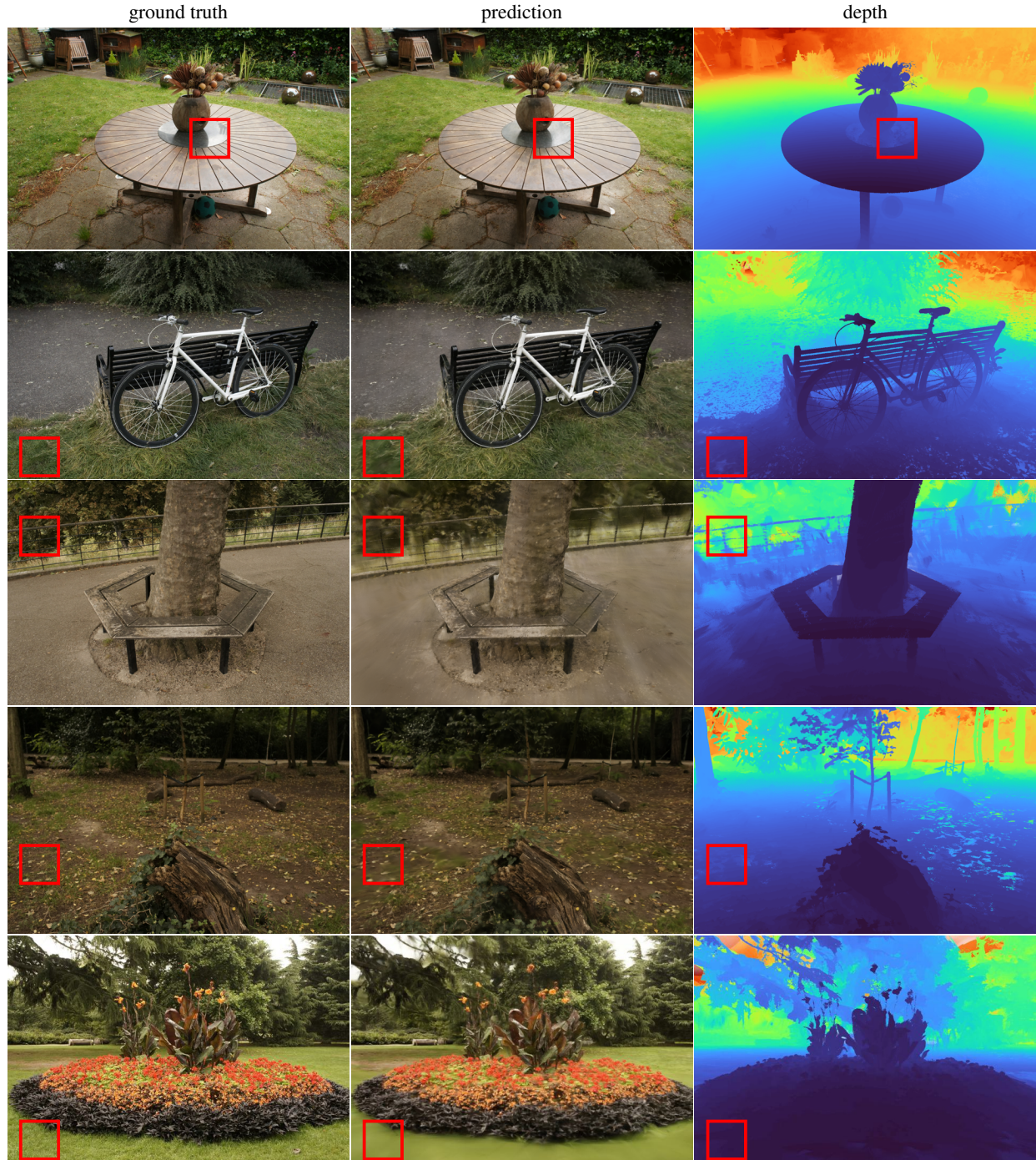


Figure 14. **Results on the outdoor scenes from mip-NeRF 360 dataset [4].** We show the **ground-truth** image, the **prediction**, and the **predicted depth map** on scenes: *garden* (top), *bicycle*, *treehill*, *stump*, and *flowers* (bottom). The red squares highlight regions where there are visible errors in the images. Tetra-NeRF is not able to represent the grass and the ground in *flowers* and *treehill* scenes well and we can see blur artefacts. On the *garden* scene, the reconstruction achieves high fidelity except for the centre of the table, where the reflections are slightly incorrect.

E. Varying the number of input points

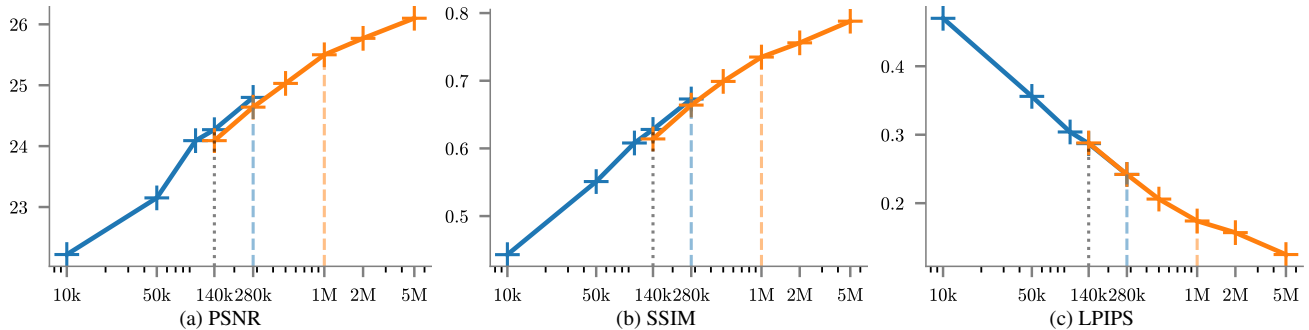


Figure 15. **Performance with different sizes and quality of input point cloud on the mip-NeRF 360 kitchen scene [4].** We show the PSNR, SSIM, and LPIPS (Alex) [75] results. The blue and orange colours are used for the sparse and dense COLMAP-reconstructed point clouds, respectively. The dashed lines show the variants used in other experiments in the paper. The dotted line shows the original number of points of the sparse point cloud (without adding any random points). We show the results for 10k, 50k, 100k, 140k (original sparse point cloud), 280k (when adding the same number of random points), 500k, 1M, 2M, and 5M points. As expected, the performance grows with the number of points.

We conducted a limited study on the effect of different sizes of the input point cloud. We consider both the sparse and dense point clouds and analyse the performance as we sub-sample the points or add more points randomly. We evaluate the results on the *garden* scene from the mip-NeRF 360 [4] dataset. However, in order to save computational resources, we only train the method for 100k iterations. The results are visualised in Figure 15. We used the sparse and dense COLMAP-reconstructed point clouds. The size of the sparse point cloud was 138,766 points, and we evaluated the following variants: 10k, 50k, and 100k randomly sub-sampled points from the full sparse point cloud, the full sparse point cloud, and the full sparse point cloud, where for each original point we added a new random point as described in the main paper. For the dense point cloud, the size was 61,766,989 points, and we sub-sampled it to the size of the sparse point cloud, twice that size, 500k, 1M, 2M, and 5M points. As expected, the performance improves with the number of points used, as it leads to a finer subdivision of the scene around the surface. For all experiments in the paper and the *Supp. Mat.*, we used random sub-sampling. An interesting direction for future work is to investigate more sophisticated strategies that, *e.g.*, sample more points around fine details such as corners and edges and fewer points in planar regions.

References

- [1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*, pages 696–712. Springer, 2020. 2, 9
- [2] Pierre Alliez, Simon Giraudot, Clément Jamin, Florent Lafarge, Quentin Mérigot, Jocelyn Meyron, Laurent Saboret, Nader Salman, Shihao Wu, and Necip Fazil Yildiran. Point set processing. In *CGAL User and Reference Manual*. CGAL Editorial Board, 5.5.1 edition, 2022. 5
- [3] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 1, 2, 4, 5, 7, 8, 9, 10, 15
- [4] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 1, 2, 5, 7, 8, 9, 15, 16, 17, 18
- [5] Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition. *arXiv preprint arXiv:2008.03824*, 2020. 2
- [6] Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. NoPe-NeRF: Optimising neural radiance field with no pose prior. *arXiv preprint arXiv:2212.07388*, 2022. 2
- [7] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. NeRD: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12684–12694, 2021. 2
- [8] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 425–432, 2001. 2

- [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensorRF: Tensorial radiance fields. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 333–350. Springer, 2022. 1, 2, 4
- [10] Anpei Chen, Zexiang Xu, Xinyue Wei, Siyu Tang, Hao Su, and Andreas Geiger. Factor fields: A unified framework for neural fields and beyond. *arXiv preprint arXiv:2302.01226*, 2023. 1, 2, 3
- [11] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. 2
- [12] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2524–2534, 2020. 2
- [13] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry- and image-based approach. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 11–20, 1996. 2
- [14] B Delaunay, S Vide, A Lamémoire, and V De Georges. Bulletin de l’academie des sciences de l’URSS. *Classe des sciences mathématiques et naturelles*, 6:793–800, 1934. 2, 3, 4, 5
- [15] Boyang Deng, Jonathan T. Barron, and Pratul P. Srinivasan. JaxNeRF: an efficient JAX implementation of NeRF, 2020. 8, 15
- [16] Michael S Floater. Generalized barycentric coordinates and applications. *Acta Numerica*, 24:161–214, 2015. 3, 4
- [17] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 1, 2, 4, 5, 6, 9, 10
- [18] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32:1362–1376, 2009. 2
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. 5
- [20] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018. 2, 8, 15
- [21] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5875–5884, 2021. 2
- [22] Vu Hoang Hiep, Renaud Keriven, Patrick Labatut, and Jean-Philippe Pons. Towards high-resolution large-scale multi-view stereo. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1430–1437. IEEE, 2009. 2
- [23] Michal Jancosek and Tomas Pajdla. Multi-view reconstruction preserving weakly-supported surfaces. In *CVPR 2011*, pages 3121–3128. IEEE, 2011. 2
- [24] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, volume 7, page 0, 2006. 2
- [25] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 5, 6, 8, 9, 13, 14
- [26] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. In *Computer Graphics Forum*, volume 40, pages 29–43. Wiley Online Library, 2021. 2, 8, 15
- [27] Jonáš Kulháněk, Erik Derner, Torsten Sattler, and Robert Babuška. ViewFormer: Nerf-free neural rendering from few images using transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 198–216. Springer, 2022. 2
- [28] Patrick Labatut, Jean-Philippe Pons, and Renaud Keriven. Efficient multi-view reconstruction of large-scale scenes using interest points, Delaunay triangulation and graph cuts. In *2007 IEEE 11th international conference on computer vision*, pages 1–8. IEEE, 2007. 2
- [29] Patrick Labatut, J-P Pons, and Renaud Keriven. Robust and efficient surface reconstruction from range data. In *Computer graphics forum*, volume 28, pages 2275–2290. Wiley Online Library, 2009. 2
- [30] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. BARF: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021. 2
- [31] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 1, 2, 3, 5, 6, 9, 13
- [32] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. In *International Conference on Learning Representations*, 2020. 5
- [33] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Transactions on Graphics (ToG)*, 38(4):1–14, 2019. 2, 6, 13
- [34] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3D surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 2

- [35] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 2
- [36] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):99–108, 1995. 1, 3
- [37] Moustafa Meshry, Dan B. Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snaveley, and Ricardo Martin-Brualla. Neural rerendering in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [38] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 15
- [39] Tomas Möller. A fast triangle-triangle intersection test. *Journal of graphics tools*, 2(2):25–30, 1997. 2
- [40] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 1, 2, 3, 4, 5, 9
- [41] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3D reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 2013. 4
- [42] Steven G Parker, James Bigler, Andreas Dietrich, Heiko Friedrich, Jared Hoberock, David Luebke, David McAllister, Morgan McGuire, Keith Morley, Austin Robison, et al. OptiX: a general purpose ray tracing engine. *ACM Transactions on Graphics (ToG)*, 29(4):1–13, 2010. 2, 4, 5
- [43] Bui Tuong Phong. Illumination for computer generated pictures. *Communications of the ACM*, 18(6):311–317, 1975. 2
- [44] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. DreamFusion: Text-to-3D using 2D diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2
- [45] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. KiloNeRF: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14335–14345, 2021. 2
- [46] Christian Reiser, Richard Szeliski, Dor Verbin, Pratul P Srinivasan, Ben Mildenhall, Andreas Geiger, Jonathan T Barron, and Peter Hedman. MERF: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. *arXiv preprint arXiv:2302.12249*, 2023. 2
- [47] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3D: Large-scale learning and evaluation of real-life 3D category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10901–10911, 2021. 2
- [48] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX 16*, pages 623–640. Springer, 2020. 2
- [49] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12216–12225, 2021. 2, 7, 8, 15
- [50] Radu Alexandru Rosu and Sven Behnke. Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [51] Darius Rückert, Linus Franke, and Marc Stamminger. ADOP: Approximate differentiable one-pixel point rendering. *ACM Transactions on Graphics (ToG)*, 41(4):1–14, 2022. 2
- [52] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 4, 5
- [53] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 4, 5
- [54] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 1, 2, 4
- [55] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretschmar. Block-NeRF: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 2
- [56] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33:7537–7547, 2020. 4
- [57] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, et al. Nerfstudio: A modular framework for neural radiance field development. *arXiv preprint arXiv:2302.04264*, 2023. 2, 5
- [58] Chengzhou Tang and Ping Tan. BA-net: Dense bundle adjustment networks. In *International Conference on Learning Representations*, 2019. 2

- [59] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-NeRF: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022. 2
- [60] Sudheendra Vijayanarasimhan, Susanna Ricco, Cordelia Schmid, Rahul Sukthankar, and Katerina Fragkiadaki. SfM-Net: Learning of structure and motion from video. *arXiv preprint arXiv:1704.07804*, 2017. 2
- [61] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-NeRF: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022. 2
- [62] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021. 2
- [63] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4690–4699, June 2021. 2
- [64] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. NeuS2: Fast learning of neural implicit surfaces for multi-view reconstruction. *arXiv preprint arXiv:2212.05231*, 2022. 2
- [65] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. NeRF-: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 2
- [66] Daniel N Wood, Daniel I Azuma, Ken Aldinger, Brian Curless, Tom Duchamp, David H Salesin, and Werner Stuetzle. Surface light fields for 3D photography. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 287–296, 2000. 2
- [67] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-NeRF: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022. 1, 2, 3, 5, 6, 9, 10, 13, 14
- [68] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. MVSNet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 2
- [69] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2
- [70] Lior Yariv, Peter Hedman, Christian Reiser, Dor Verbin, Pratul P Srinivasan, Richard Szeliski, Jonathan T Barron, and Ben Mildenhall. BakedSDF: Meshing neural SDFs for real-time view synthesis. *arXiv preprint arXiv:2302.14859*, 2023. 2
- [71] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 1, 2
- [72] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. PixelNeRF: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 2
- [73] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. MonoSDF: Exploring monocular geometric cues for neural implicit surface reconstruction. In *Advances in Neural Information Processing Systems*, 2022. 2
- [74] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. NeRF++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2, 8, 15
- [75] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 6, 8, 10, 13, 15, 18
- [76] Chaoqiang Zhao, Qiyu Sun, Chongzhen Zhang, Yang Tang, and Feng Qian. Monocular depth estimation based on deep learning: An overview. *Science China Technological Sciences*, 63(9):1612–1627, 2020. 3