

# Behind the Scenes: Density Fields for Single View Reconstruction

Felix Wimbauer<sup>1,2</sup>

Nan Yang<sup>1</sup>

Christian Rupprecht<sup>3</sup>

Daniel Cremers<sup>1,2,3</sup>

<sup>1</sup>Technical University of Munich

<sup>2</sup>MCML

<sup>3</sup>University of Oxford

{felix.wimbauer, nan.yang, cremers}@tum.de

chrisr@robots.ox.ac.uk

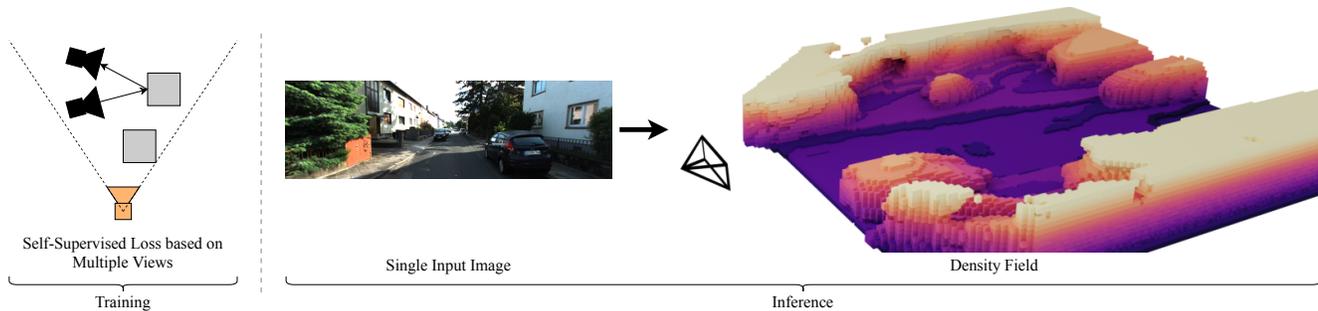


Figure 1. **Predicting a Density Field from a Single Image.** Through a novel “density field” formulation, which decouples geometry from color, architectural improvements, and a novel self-supervised training scheme, our method learns to predict a volumetric scene representation from a single image in challenging conditions. The voxel occupancy view shows that our method predicts accurate density even in occluded regions, which is not possible in traditional depth prediction. Please check out our project page at [fwmb.github.io/bts/](https://fwmb.github.io/bts/).

## Abstract

*Inferring a meaningful geometric scene representation from a single image is a fundamental problem in computer vision. Approaches based on traditional depth map prediction can only reason about areas that are visible in the image. Currently, neural radiance fields (NeRFs) can capture true 3D including color, but are too complex to be generated from a single image. As an alternative, we propose to predict an implicit density field from a single image. It maps every location in the frustum of the image to volumetric density. By directly sampling color from the available views instead of storing color in the density field, our scene representation becomes significantly less complex compared to NeRFs, and a neural network can predict it in a single forward pass. The network is trained through self-supervision from only video data. Our formulation allows volume rendering to perform both depth prediction and novel view synthesis. Through experiments, we show that our method is able to predict meaningful geometry for regions that are occluded in the input image. Additionally, we demonstrate the potential of our approach on three datasets for depth prediction and novel-view synthesis.*

## 1. Introduction

The ability to infer information about the geometric structure of a scene from a single image is of high importance for a wide range of applications from robotics to aug-

mented reality. While traditional computer vision mainly focused on reconstruction from multiple images, in the deep learning age the challenge of inferring a 3D scene from merely a single image has received renewed attention.

Traditionally, this problem has been formulated as the task of predicting per-pixel depth values (*i.e.* depth maps). One of the most influential lines of work showed that it is possible to train neural networks for accurate single-image depth prediction in a self-supervised way only from video sequences. [14–16, 29, 44, 51, 58, 59, 61] Despite these advances, depth prediction methods are *not* modeling the true 3D of the scene: they model only a *single* depth value per pixel. As a result, it is not directly possible to obtain depth values from views other than the input view without considering interpolation and occlusion. Further, the predicted geometric representation of the scenes does not allow reasoning about areas that lie *behind* another object in the image (*e.g.* a house behind a tree), inhibiting the applicability of monocular depth estimation to 3D understanding.

Due to the recent advance of 3D neural fields, the related task of novel view synthesis has also seen a lot of progress. Instead of directly reasoning about the scene geometry, the goal here is to infer a representation that allows rendering views of the scene from novel viewpoints. While geometric properties can often be inferred from the representation, they are usually only a side product and lack visual quality.

Even though neural radiance field [32] based methods achieve impressive results, they require many training im-

ages per scene and do not generalize to new scenes. To enable generalization, efforts have been made to condition the neural network on global or local scene features. However, this has only been shown to work well on simple scenes, for example, scenes containing an object from a single category [43, 57]. Nevertheless, obtaining a neural radiance field from a single image has not been achieved before.

In this work, we tackle the problem of inferring a geometric representation from a single image by generalizing the depth prediction formulation to a continuous density field. Concretely, our architecture contains an encoder-decoder network that predicts a dense feature map from the input image. This feature map locally conditions a density field inside the camera frustum, which can be evaluated at any spatial point through a multi-layer perceptron (MLP). The MLP is fed with the coordinates of the point and the feature sampled from the predicted feature map by reprojecting points into the camera view. To train our method, we rely on simple image reconstruction losses.

Our method achieves robust generalization and accurate geometry prediction even in very challenging outdoor scenes through three key novelties:

**1. Color sampling.** When performing volume rendering, we sample color values directly from the input frames through reprojection instead of using the MLP to predict color values. We find that only predicting density drastically reduces the complexity of the function the network has to learn. Further, it forces the model to adhere to the multi-view consistency assumption during training, leading to more accurate geometry predictions.

**2. Shifting capacity to the feature extractor.** In many previous works, an *encoder* extracts image features to condition local appearance, while a high-capacity MLP is expected to generalize to multiple scenes. However, on complex and diverse datasets, the training signal is too noisy for the MLP to learn meaningful priors. To enable robust training, we significantly reduce the capacity of the MLP and use a more powerful *encoder-decoder* that can capture the entire scene in the extracted features. The MLP then only evaluates those features locally.

**3. Behind the Scenes loss formulation.** The continuous nature of density fields and color sampling allow us to reconstruct a novel view from the colors of any frame, not just the input frame. By applying a reconstruction loss between two frames that both observe areas occluded in the input frame, we train our model to predict meaningful geometry *everywhere* in the camera frustum, not just the visible areas.

We demonstrate the potential of our new approach in a number of experiments on different datasets regarding the aspects of capturing true 3D, depth estimation, and novel view synthesis. On KITTI [12] and KITTI-360 [26], we show both qualitatively and quantitatively that our model can indeed capture true 3D, and that our model

achieves state-of-the-art depth estimation accuracy. On RealEstate10K [45] and KITTI, we achieve competitive novel view synthesis results, even though our method is purely geometry-based. Further, we perform thorough ablation studies to highlight the impact of our design choices.

## 2. Related Work

In the following, we review the most relevant works that are related to our proposed method.

### 2.1. Single-Image Depth Prediction

One of the predominant formulations to capture the geometric structure of a scene from a single image is predicting a per-pixel depth map. Learning-based methods have proven able to overcome the inherent ambiguities of this task by correlating contextual cues extracted from the image with certain depth values. One of the most common ways to train a method for single-image depth prediction is to immediately regress the per-pixel ground-truth depth values [10, 27]. Later approaches supplemented the fully-supervised training with reconstruction losses [21, 56], or specialise the architecture and loss formulation [1, 11, 22, 24, 25, 54]. To overcome the need for ground-truth depth annotations, several papers focused on relying exclusively on reconstruction losses to train prediction networks. Both temporal video frames [61] and stereo frames [13], as well as combinations of both [14, 59] can be used as the reconstruction target. Different follow-up works refine the architecture and loss [15, 16, 29, 44, 51, 58]. [60] first predicts a discrete density volume as an intermediate step, from which depth maps can be rendered from different views. While they use this density volume for regularization, their focus is on improving depth prediction and their method does not demonstrate the ability to learn true 3D.

### 2.2. Neural Radiance Fields

Many works have investigated alternative approaches to representing scenes captured from a single or multiple images, oftentimes with the goal of novel view synthesis. Recently, [32] proposed to represent scenes as neural radiance fields (NeRFs). In NeRFs, a multi-layer perceptron (MLP) is optimized per scene to map spatial coordinates to color (appearance) and density (geometry) values. By evaluating the optimized MLP along rays and then integrating the color over the densities, novel views can be rendered under the volume rendering formulation [30]. Training data consists of a large number of images of the same scene from different viewpoints with poses computed from traditional SFM and SLAM methods [4, 40, 41]. The training goal is to reconstruct these images as accurately as possible. NeRF’s impressive performance inspired many follow-up works, which improve different parts of the architecture [2, 3, 7, 18, 20, 35, 38].

In the traditional NeRF formulation, an entire scene is captured in a single, large MLP. Thus, the trained network cannot be adapted to different settings or used for other scenes. Further, the MLP has to have a high capacity, resulting in slow inference. Several methods propose to condition such MLPs on feature grids or voxels [6, 28, 31, 34, 37, 43, 47, 57]. Through this, the MLP needs to store less information and can be simplified, speeding up inference [6, 28, 34, 47]. Additionally, this allows for some generalization to new scenes [33, 43, 57]. However, generalization is mostly limited to a single object category, or simple synthetic data, where the scenes differ in local details. In contrast, our proposed method can generalize to highly complex outdoor scenes. [5] also improves generalization through depth supervision and improved ray sampling.

### 2.3. Single Image Novel View Synthesis

While traditional NeRF-based methods achieve impressive performance when provided with enough images per scene, they do not work with only a single image of a scene available. In recent years, a number of methods for novel-view synthesis (NVS) from a single image emerged.

Several methods [8, 9, 49] predict layered depth images (LDI) [42] for rendering. Later approaches [46, 48] directly produce a multiplane image (MPI) [62]. [23] predicts a generalized multiplane image. Instead of directly outputting the discrete layers, the architecture’s decoder receives a variable depth value, for which it outputs the layer. In [52], a network predicts both a per-pixel depth and feature map, which are then used in a neural rendering framework. Other works [53, 55] perform image decomposition, followed by classical rendering. While these methods achieve impressive NVS results, the quality of predicted geometry usually falls short. Some methods even predict novel views without any geometric representation [39, 63].

## 3. Method

In the following, we describe a neural network architecture that predicts the geometric structure of a scene from a single image  $\mathbf{I}_1$ , as shown in Fig. 2. We first cover how we represent a scene as a continuous density field, and then propose a training scheme that allows our architecture to learn geometry even in occluded areas.

### 3.1. Notation

Let  $\mathbf{I}_1 \in [0, 1]^{3 \times H \times W} = (\mathbb{R}^3)^\Omega$  be the input image, defined on a lattice  $\Omega = \{1, \dots, H\} \times \{1, \dots, W\}$ .  $T_1 \in \mathbb{R}^{4 \times 4}$  and  $K_1 \in \mathbb{R}^{3 \times 4}$  are the corresponding world-to-camera pose matrix and projection matrix, respectively. During training, we have available an additional set of  $N = \{1, 2, \dots, n\}$  frames  $\mathbf{I}_k, k \in N$  with corresponding world-to-camera pose and projection matrices  $T_k, K_k, k \in N$ . When assuming homogeneous coordinates, a point  $\mathbf{x} \in \mathbb{R}^3$

in world coordinates can be projected onto the image plane of frame  $k$  with the following operation:  $\pi_k(\mathbf{x}) = K_k T_k \mathbf{x}$

### 3.2. Predicting a Density Field

We represent the geometric structure of a scene as a function, which maps scene coordinates  $\mathbf{x}$  to volume density  $\sigma$ . We term this function “density field”. Inference happens in two steps. From the input image  $\mathbf{I}_1$ , an encoder-decoder network first predicts a pixel-aligned feature map  $\mathbf{F} \in (\mathbb{R}^C)^\Omega$ . The idea behind this is that every feature  $f_{\mathbf{u}} = \mathbf{F}(\mathbf{u})$  at pixel location  $\mathbf{u} \in \Omega$  captures the distribution of local geometry along the ray from the camera origin through the pixel at  $\mathbf{u}$ . It also means that the density field is designed to lie inside the camera frustum. For points outside of this frustum, we extrapolate features from within the frustum.

To obtain a density value at a 3D coordinate  $\mathbf{x}$ , we first project  $\mathbf{x}$  onto the input image  $\mathbf{u}'_1 = \pi_1(\mathbf{x})$  and bilinearly sample the feature  $f_{\mathbf{u}'_1} = \mathbf{F}(\mathbf{u}'_1)$  at that position. This feature  $f_{\mathbf{u}'_1}$ , along with the positional encoding [32]  $\gamma(d)$  of the distance  $d$  between  $\mathbf{x}$  and the camera origin, and the positional encoding  $\gamma(\mathbf{u}'_1)$  of the pixel, is then passed to a multi-layer perceptron (MLP)  $\phi$ . During training,  $\phi$  and  $\mathbf{F}$  learn to describe the density of the scene given the input view. We can interpret the feature representation  $f_{\mathbf{u}'_1}$  as a descriptor of the density along a ray through the camera center and pixel  $\mathbf{u}'_1$ . In turn,  $\phi$  acts as a decoder, that given  $f_{\mathbf{u}'_1}$  and a distance to the camera, predicts the density at the 3D location  $\mathbf{x}$ .

$$\sigma_{\mathbf{x}} = \phi(f_{\mathbf{u}'_1}, \gamma(d), \gamma(\mathbf{u}'_1)) \quad (1)$$

Unlike most current works on neural fields, we *do not* use  $\phi$  to also predict color. This drastically reduces the complexity of the distribution along a ray as density distributions tend to be simple, while color often contains complex high-frequency components. In our experiments, this makes capturing such a distribution in a single feature, so that it can be evaluated by an MLP, much more tractable.

### 3.3. Volume Rendering with Color Sampling

When rendering the scene from a novel viewpoint, we do not retrieve color from our scene representation directly. Instead, we sample the color for a point in 3D space from the available images. Concretely, we first project a point  $\mathbf{x}$  into a frame  $k$  and then bilinearly sample the color  $c_{\mathbf{x},k} = \mathbf{I}_k(\pi_k(\mathbf{x}))$ .

By combining  $\sigma_{\mathbf{x}}$  and  $c_{\mathbf{x},k}$ , we can perform volume rendering [19, 30] to synthesize novel views. We follow the discretization strategy of other radiance field-based methods, e.g. [32]. To obtain the color  $\hat{c}_k$  for a pixel in a novel view, we emit a ray from the camera and integrate the color along the ray over the probability of the ray ending at a certain distance. To approximate this integral, density and color are evaluated at  $S$  discrete steps  $\mathbf{x}_i$  along the ray. Let  $\delta_i$  be the distance between  $\mathbf{x}_i$  and  $\mathbf{x}_{i+1}$ , and  $\alpha_i$  be the probability of

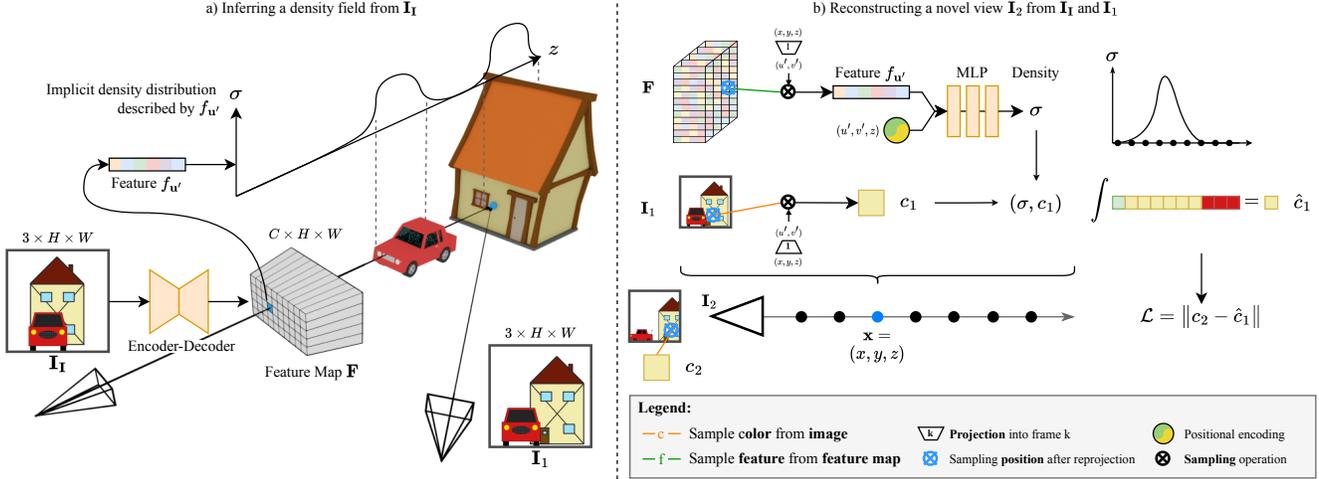


Figure 2. **Overview.** *a)* Our method first predicts a pixel-aligned feature map  $\mathbf{F}$ , which describes a density field, from the input image  $\mathbf{I}_1$ . For every pixel  $\mathbf{u}'$ , the feature  $f_{\mathbf{u}'}$  implicitly describes the density distribution along the ray from the camera origin through  $\mathbf{u}'$ . Crucially, this distribution can model density even in occluded regions (*e.g.* the house). *b)* To render novel views, we perform volume rendering. For any point  $\mathbf{x}$ , we project  $\mathbf{x}$  into  $\mathbf{F}$  and sample  $f_{\mathbf{u}'}$ . This feature is combined with positional encoding and fed into an MLP to obtain density  $\sigma$ . We obtain the color  $c$  by projecting  $\mathbf{x}$  into one of the views, in this case,  $\mathbf{I}_1$ , and directly sampling the image.

a ray ending between  $\mathbf{x}_i$  and  $\mathbf{x}_{i+1}$ . From the previous  $\alpha_j$ s, we can compute the probability  $T_i$  that the ray does not terminate before  $\mathbf{x}_i$ , *i.e.* the probability that  $\mathbf{x}_i$  is not occluded.

$$\alpha_i = \exp(1 - \sigma_{\mathbf{x}_i} \delta_i) \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (2)$$

$$\hat{c}_k = \sum_{i=1}^S T_i \alpha_i c_{\mathbf{x}_i, k} \quad (3)$$

Similarly, we can also retrieve the expected ray termination depth, which corresponds to the depth in a depth map. Let  $d_i$  be the distance between  $\mathbf{x}_i$  and the ray origin.

$$\hat{d} = \sum_{i=1}^S T_i \alpha_i d_i \quad (4)$$

This rendering formulation is very flexible. We can sample the color values from any frame, and, crucially, it can be a different frame from the input frame. It is even possible to obtain multiple colors from multiple different frames for a single ray, which enables reasoning about occluded areas during training. Note that even though different frames can be used, the density is always based on features from the input image and does not change. During inference, color sampling from different frames is not necessary, everything can be done based on a single input image.

### 3.4. Behind the Scenes Loss Formulation

Our training goal is to optimize both the encoder-decoder network and  $\phi$  to predict a density field only from the input image, such that it allows the reconstruction of other views.

Similar to radiance fields and self-supervised depth prediction methods, we rely on an image reconstruction loss. For a single sample, we first compute the feature map  $\mathbf{F}$  from  $\mathbf{I}_1$  and randomly partition *all* frames  $\hat{N} = \{\mathbf{I}_1\} \cup N$  into two sets  $N_{\text{loss}}, N_{\text{render}}$ . Note that the input image can end up in any of the two sets. We reconstruct the frames in  $N_{\text{loss}}$  by sampling colors from  $N_{\text{render}}$  using the camera poses and the predicted densities. The photometric consistency between the reconstructed frames and the frames in  $N_{\text{loss}}$  serves as supervision for the density field. In practice, we randomly sample  $p$  patches  $P_i$  to use patch-wise photometric measurement. For every patch  $P_i$  in  $N_{\text{loss}}$ , we obtain a reconstructed patch  $\hat{P}_{i,k}$  from *every* frame  $k \in N_{\text{render}}$ . We aggregate the costs between  $P_i$  and every  $\hat{P}_{i,k}$  by taking the per-pixel *minimum* across the different frames  $k$ , similar to [14]. The intuition behind this is that for every patch, there is a frame in  $N_{\text{render}}$ , which “sees” the same surface. Therefore, if the predicted density is correct, then it results in a very good reconstruction and a low error.

For the final loss formula, we use a combination of L1 and SSIM [50] to compute the photometric discrepancy, as well as an edge-aware smoothness term. Let  $d_i^*$  denote the inverse, mean-normalized expected ray termination depth of patch  $P_i$ . Both  $\mathcal{L}_{\text{ph}}$  and  $\mathcal{L}_{\text{eas}}$  are computed per  $(x, y)$  element of the patch, thus resulting in 2D loss maps. They are then aggregated when computing  $\mathcal{L}$ .

$$\mathcal{L}_{\text{ph}} = \min_{k \in N_{\text{render}}} \left( \lambda_{\text{L1}} \text{L1}(P_i, \hat{P}_{i,k}) + \lambda_{\text{SSIM}} \text{SSIM}(P_i, \hat{P}_{i,k}) \right) \quad (5)$$

$$\mathcal{L}_{\text{eas}} = |\delta_x d_i^*| e^{-|\delta_x P_i|} + |\delta_y d_i^*| e^{-|\delta_y P_i|} \quad (6)$$

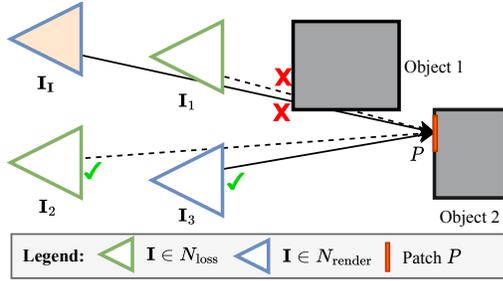


Figure 3. **Loss in Ocluded Regions.** Patch  $P$  on Object 2 is occluded by Object 1 in the input frame  $\mathbf{I}_1$  and  $\mathbf{I}_1$ . In order to correctly reconstruct  $P$  in  $\mathbf{I}_2$  from  $\mathbf{I}_3$ , the network needs to predict density for Object 2 *behind* Object 1.

$$\mathcal{L} = \sum_{i=1}^P \sum_{x,y \in P} (\mathcal{L}_{\text{ph}} + \lambda_{\text{eas}} \mathcal{L}_{\text{eas}})(x, y) \quad (7)$$

**Learning true 3D.** Our loss formula Eq. (7) is the same as for self-supervised depth prediction methods, like [14]. The key difference, however, is that depth prediction methods can only densely reconstruct the input image, for which the per-pixel depth was predicted.

In contrast, our density field formulation allows us to reconstruct *any* frame from *any other* frame. Consider an area of the scene, which is occluded in the input  $\mathbf{I}_1$ , but visible in two other frames  $\mathbf{I}_k, \mathbf{I}_{k+1}$ , as depicted in Fig. 3: During training, we aim to reconstruct this area in  $\mathbf{I}_k$ . The reconstruction based on colors sampled from  $\mathbf{I}_{k+1}$  will give a clear training signal to correctly predict the geometric structure of this area, even though it is occluded in  $\mathbf{I}_1$ . Note, that in order to learn geometry about occluded areas, we require at least **two additional** views besides the input during training, *i.e.* to look *behind the scenes*.

**Handling invalid samples.** While the frustums of the different views overlap for the most part, there is still a chance of a ray leaving the frustums, thus sampling invalid features, or sampling invalid colors. Such invalid rays lead to noise and instability in the training process. Therefore, we propose a policy to detect and remove invalid rays. Our intuition is that when the amount of contribution to the final aggregated color, that comes from invalidly sampled colors or features, exceeds a certain threshold  $\tau$ , the ray should be discarded. Consider a ray that is evaluated at positions  $\mathbf{x}_i, i \in [1, 2, \dots, S]$  and reconstructed from frames  $K: O_{i,k}, k \in \{\mathbf{I}\} \cup K$  denotes the indicator function that  $\mathbf{x}_i$  is outside the camera frustum of frame  $k$ . Note that we always sample features from the input frame. We define  $\text{IV}(k)$  to be the function indicating that the rendered color based on frame  $k$  is invalid as:

$$\text{IV}(k) = \sum_{i=1}^S T_i \alpha_i (O_{i,1} \vee O_{i,k}) > \tau \quad (8)$$

Only if  $\text{IV}(k)$  is true for *all* frames the ray was reconstructed from, we ignore the ray when computing the loss value. The reasoning behind this is that non-invalid rays will still lead to the lowest error. Therefore, the min operation in Eq. (5) will ignore the invalid rays.

### 3.5. Implementation Details

We implement our model in PyTorch [36] on a single Nvidia RTX A40 GPU with 48GB memory. The encoder-decoder network follows [14] using a ResNet encoder [17] and predicts feature maps with 64 channels. The MLP  $\phi$  is made lightweight with only 2 fully connected layers and 64 hidden nodes each. We use a batch size of 16 and sample 32 patches of size  $8 \times 8$  from the images for which we want to compute the reconstruction loss. Every ray is sampled at 64 locations, based on linear spacing in inverse depth. For more details, *e.g.* exact network architecture and further hyperparameters, please refer to the supplementary material.

## 4. Experiments

To demonstrate the abilities and advantages of our proposed method, we conduct a wide range of experiments. First, we demonstrate that our method is uniquely able to capture a holistic geometry representation of the scene, even in areas that are occluded in the input image. Additionally, we also show the effect of different data setups on the prediction quality. Second, we show that our method, even though depth maps are only a side product of our scene representation, achieves depth accuracy on par with other state-of-the-art self-supervised methods, that are specifically designed for depth prediction. Third, we demonstrate that, even though our representation is geometry-only, our method can be used to perform high-quality novel view synthesis from a single image. Finally, we conduct thorough ablation studies based on occupancy estimation and depth prediction to justify our design choices.

### 4.1. Data

For our experiments, we use three different datasets: KITTI [12], KITTI-360 [26] (autonomous driving), and RealEstate10K [62] (indoor). RealEstate10K only has monocular sequences, while KITTI and KITTI-360 provide stereo. KITTI-360 also contains fisheye camera frames facing left and right. For monocular data, we use three timesteps per sample, for stereo sequences (possibly with fisheye frames), we use two timesteps. The fisheye frames are offset by one second to increase the overlap of the different camera frustums.<sup>1</sup> Training is performed for 50 epochs on KITTI (approx. 125k steps), 25 epochs on KITTI-360 (approx. 150k steps), and 360k iterations on RealEstate10K.

<sup>1</sup>More details on offsets, pose data, and data splits in the supp. mat.

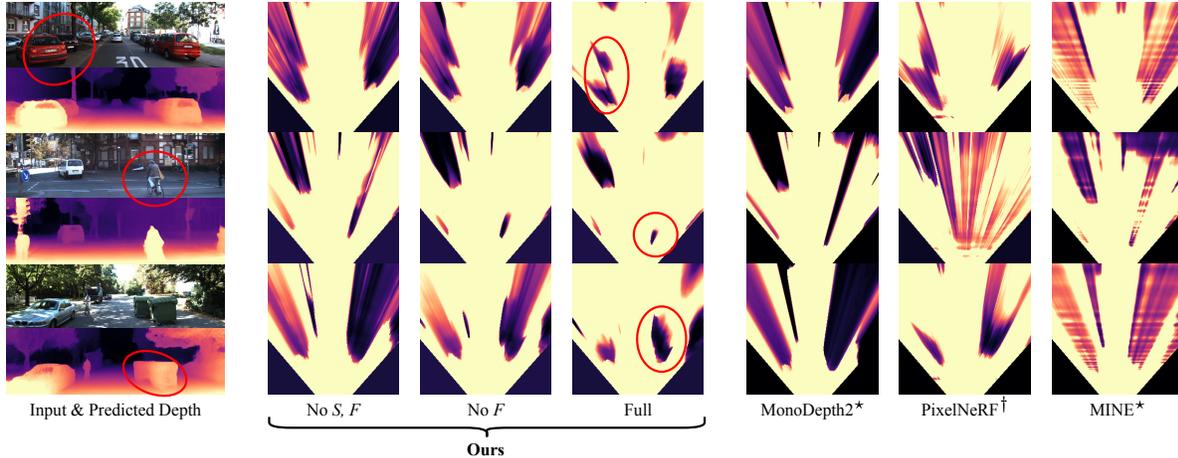


Figure 4. **Occupancy Estimation.** Top-down visualization of predicted occupancy volumes. We show an area of  $x = [-9m, 9m]$ ,  $z = [3m, 21m]$  and  $y = [0m, 1m]$  (just above the road). Our method produces an accurate volumetric reconstruction, even for occluded regions. Training with more views improves the quality. Depth prediction methods like MonoDepth2 [14] do not predict a full 3D volume. Thus, objects cast “occupancy shadows” behind them. Volumetric methods like PixelNeRF [57] and MINE [23] produce noisy predictions. Inference is from a single image. **Legend:** *S*: Stereo, *F*: Fisheye,  $\star$ : official checkpoint,  $\dagger$ : trained in same setup as our Full variant.

Method	$O_{acc} \uparrow$	$IE_{acc} \uparrow$	$IE_{rec} \uparrow$
Depth $^\dagger$ [14]	<b>0.94</b>	n/a	n/a
Depth $^\dagger$ + 4m [14]	0.91	0.63	<u>0.22</u>
PixelNeRF $^\dagger$ [57]	<u>0.92</u>	0.63	<b>0.43</b>
<b>Ours</b> (No <i>S</i> , <i>F</i> )	<b>0.94</b>	0.70	0.06
<b>Ours</b> (No <i>F</i> )	<b>0.94</b>	<u>0.71</u>	0.09
<b>Ours</b>	<b>0.94</b>	<b>0.77</b>	<b>0.43</b>

Table 1. **3D Scene Occupancy Accuracy on KITTI-360.** We evaluate the capability of the model to predict occupancy *behind* objects in the image. Ground truth occupancy maps are computed from 20 consecutive Lidar scans per frame. Depth prediction [14] naturally has no ability to predict behind occlusions. PixelNeRF [57] can predict free space in occluded regions, but produces poor overall geometry. Our method improves when training with more views. Inference from a single image. Samples are evenly spaced in a cuboid  $w = [-4m, 4m]$ ,  $h = [-1m, 0m]$ ,  $d = [3m, 20m]$  relative to the camera. **Legend:** ref. Fig. 4.

We use a resolution of  $640 \times 192$  for KITTI and KITTI-360, and follow [23] in using a resolution of  $384 \times 256$  for RealEstate10K.

## 4.2. Capturing true 3D

Evaluation of fully geometric 3D representations like density fields is difficult. Real-world datasets usually only provide ground truth data captured from a single viewpoint, *e.g.* RGB-D frames and Lidar measurements. Nevertheless, we aim to evaluate and compare this key advantage of our method both qualitatively and quantitatively. Through our proposed training scheme, our networks are able to learn to also predict meaningful geometry in occluded areas.

To overcome the lack of volumetric ground truth, we ac-

cumulate Lidar scans to build reference occupancy maps for KITTI-360. Consider a single input frame for which we want to evaluate an occupancy prediction: As KITTI-360 is a driving dataset with a forward-moving camera, the consecutive Lidar scans captured a short time later measure different areas within the camera frustum. Note that these Lidar measurements can reach areas that are occluded in the input image. To determine whether a point is occupied, we check whether it is *in front* of the measured surface for any of the Lidar scans. Intuitively, every Lidar measurement “carves out” unoccupied areas in 3D space. By accumulating enough Lidar scans, we obtain a reliable occupancy measurement of the entire camera frustum. Whether a point is visible in the input frame can be checked using the Lidar scan corresponding to the input frame.<sup>2</sup>

For every frame, we sample points in a cuboid area in the camera frustum and compute the following metrics: 1. Occupancy accuracy ( $O_{Acc}$ ), 2. Invisible and empty accuracy ( $IE_{Acc}$ ), and 3. Invisible and empty recall ( $IE_{Rec}$ ).  $O_{Acc}$  evaluates the occupancy predictions across the whole scene volume.  $IE_{Acc}$  and  $IE_{Rec}$  specifically evaluate invisible regions, evaluating performance beyond depth prediction.

We train a MonoDepth2 [14] model to serve as a baseline representing ordinary depth prediction methods. Here, we consider all points behind the predicted depth to be occupied. Additionally, we evaluate a version in which we consider points only up to 4 meters (average car length) behind the predicted depth as occupied. As a second baseline, we train a PixelNeRF [57] model, one of the most prominent NeRF variants that also has the ability to generalize.

To demonstrate that our loss formulation generates

<sup>2</sup>More details on the exact procedure and examples in the supp. mat.

Model	Volum.	Split	Abs Rel ↓	RMSE ↓	$\alpha < 1.25$ ↑
PixelNeRF [57]	✓		0.130	5.134	0.845
EPC++ [29]	✗		0.128	5.585	0.831
MonoDepth2 [14]	✗		0.106	4.750	0.874
PackNet [16]	✗	Eigen [10]	0.111	4.601	0.878
DepthHint [51]	✗		0.105	4.627	0.875
FeatDepth [44]	✗		0.099	4.427	0.889
DevNet [60]	(✓)		<b>0.095</b>	<b>4.365</b>	<b>0.895</b>
<b>Ours</b>	✓		0.102	<u>4.407</u>	0.882
MINE [23]	✓	Tuls. [49]	0.137	6.592	0.839
<b>Ours</b>	✓		<b>0.132</b>	<b>6.104</b>	<b>0.873</b>

Table 2. **Depth Prediction on KITTI.** While our goal is full volumetric scene understanding, we compare to state-of-the-art self-supervised depth estimation method. Our approach achieves competitive performance while clearly improving over other volumetric approaches like PixelNeRF [57] and MINE [23]. DevNet [60] performs better, but does not show any results of their volume.

strong training signals for occluded regions, given the right data, we train our model in several different data configurations. By removing the fisheye, respectively fisheye, and stereo frames, the training signal for occluded areas becomes much weaker. Tab. 1 reports the obtained results.

The depth prediction baselines achieve a strong overall accuracy, but are, by design, not able to predict meaningful free space in occluded areas. PixelNeRF can predict free space in occluded areas but produces poor overall geometry. Our model achieves strong overall accuracy, while it is also able to recover the geometry of the occluded parts of the scene. Importantly, our model becomes better at predicting *free space in occluded areas* when training with more views, naturally providing a better training signal for occluded areas. To qualitatively visualize these results we sample the camera frustum in horizontal slices from the center of the image downwards and aggregate the density in Fig. 4. This shows the layout of the scene, similar to the birds-eye perspective but for density. In the Full variant, the strong signal lets our model learn sharp object boundaries, as can be seen for several cars in the examples. For depth prediction, all objects cast occupancy shadows along the viewing direction. PixelNeRF predicts a volumetric representation with free space in occluded regions. However, the results are noisy and the geometry is inaccurate. MINE [23] also specializes in predicting a volumetric representation from a single image. However, it does not produce meaningful density prediction behind objects. Instead, similar to depth prediction, all objects cast occupancy shadows along the viewing direction.

### 4.3. Depth Prediction

While our method does not predict depth maps directly, they can be synthesized as a side product from our representation through the expected ray termination depth  $\hat{d}$ . To

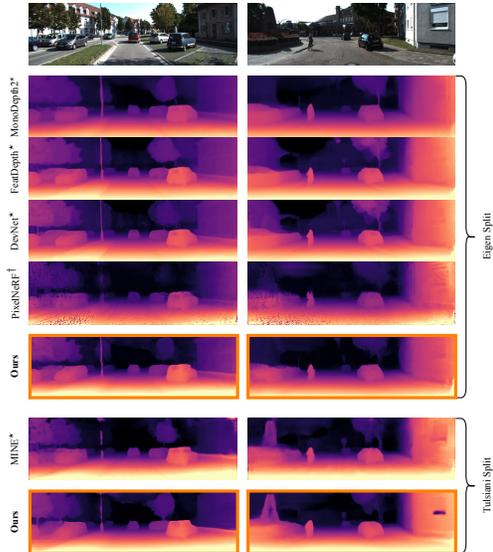


Figure 5. **Depth Prediction on KITTI.** Expected ray termination depth compared with depth prediction results of other state-of-the-art methods [14, 23, 44, 57, 60] on both the Eigen [10] and [49] split. Our predictions are very detailed and sharp, and capture the structure of the scene, even when trained on a smaller split such as Tulsiani. Visualizations for DevNet and FeatDepth are taken from [60]. **Legend:** ref. Fig. 4.

demonstrate that our predicted representation achieves high accuracy, we train our model on KITTI sequences and compare to both self-supervised depth prediction methods and volume reconstruction methods.

As can be seen in Tab. 2 and Fig. 5, our method performs on par with the current state-of-the-art methods for self-supervised depth prediction. Our synthesized depth maps capture finer details and contain fewer artifacts, as often seen with depth maps obtained from neural radiance field-based methods, like PixelNeRF [57] and MINE [23]. Overall, we achieve competitive performance, even though depth prediction is not the main objective of our approach.

### 4.4. Novel View Synthesis from a Single Image

As we obtain a volumetric representation of a scene from a single image, we are able to synthesize images from novel viewpoints by sampling color from the input image. Thus, we also evaluate novel view synthesis from a single image. To demonstrate the variability of our approach, we train two models, one on RealEstate10K [62], and one on the KITTI (Tulsiani split [49]). As Tab. 4 shows, our method achieves strong performance on both datasets, despite the fact, that we only predict geometry and obtain color by sampling the input image. Our results are comparable with many recent methods, that were specifically designed for this task, and of which some even use sparse depth supervision during training for RealEstate10K (MPI, MINE). MINE [23] achieves

Method	Configuration			Occupancy Estimation			Depth Prediction		
	Features	MLP	Predicts	O <sub>acc</sub> ↑	IE <sub>acc</sub> ↑	IE <sub>rec</sub> ↑	Abs Rel ↓	RMSE ↓	$\alpha < 1.25$ ↑
Baselines	Enc	Big	$\sigma + c$	0.92	0.63	<u>0.41</u>	0.130	5.134	0.845
	E+D	Big	$\sigma + c$	0.93	0.62	<b>0.43</b>	0.149	5.441	0.800
	Enc	Small	$\sigma + c$	0.92	<u>0.69</u>	0.31	0.112	4.897	0.860
	E+D	Small	$\sigma + c$	0.93	<u>0.69</u>	0.15	0.109	4.758	0.864
	Enc	Small	$\sigma$	<b>0.94</b>	<b>0.77</b>	0.39	<u>0.105</u>	4.590	<u>0.872</u>
<b>Ours</b>	<b>E+D</b>	<b>Small</b>	$\sigma$	<b>0.94</b>	<b>0.77</b>	<b>0.43</b>	<b>0.102</b>	<b>4.407</b>	<b>0.882</b>
<b>Ours</b>	Keep invalid rays			0.94	0.77	0.41	0.108	4.493	0.875

Table 3. **Ablation Studies.** Evaluation of variants with different contributions (predicting only density  $\sigma$  and sampling color, shifting capacity from the MLP to the feature extractor, discarding invalid rays) turned on / off. Occupancy estimation results on KITTI-360 and depth prediction results on KITTI. The variant using only an encoder, big MLP, and color prediction corresponds exactly to the PixelNeRF [57] architecture, but with our training scheme. **Legend:** *Enc* Encoder, *E+D* Encoder-Decoder,  $\sigma$  density,  $c$  color.

Model	KITTI			RealEstate10K		
	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
SynSin [52]	n/a	n/a	n/a	1.180	0.740	22.3
Tulsiani [49]	n/a	0.572	16.5	<u>0.176</u>	<u>0.785</u>	23.5
MPI [48]	n/a	0.733	19.5	n/a	n/a	n/a
MINE [23]	<b>0.112</b>	<b>0.828</b>	<b>21.9</b>	<b>0.156</b>	<b>0.822</b>	<b>24.5</b>
PixelNeRF [57]	0.175	0.761	<u>20.1</u>	n/a	n/a	n/a
<b>Ours</b>	<u>0.144</u>	<u>0.764</u>	<u>20.1</u>	0.194	0.755	<u>24.0</u>

Table 4. **Novel View Synthesis.** We test the NVS ability on KITTI (Tulsiani split [49]) and RealEstate10K (MINE split [23], target frame randomly sampled within 30 frames). Even though our method does not predict color, we still achieve strong results.

slightly better accuracy. This can be attributed to them being able to predict color and thereby circumventing issues arising from imperfect geometry.

## 4.5. Ablation Studies

Our architectural design choices are critically important for the strong performance of our method. To quantify the impact of the different contributions, we conduct ablation studies based on occupancy estimation on KITTI-360 and depth prediction on KITTI. PixelNeRF [57] can be seen as a basis, which we modify step-by-step to reach our proposed model. Namely, we 1. shift capacity from the MLP to the feature extractor and 2. introduce color sampling as an alternative to predicting the color alongside density.

As Tab. 3 shows, reducing the MLP capacity and using a more powerful encoder-decoder rather than encoder as a feature extractor allows the model to learn significantly more precise overall geometry. We conjecture that a powerful feature extractor is more suited to generalize to unseen scenes based on a single input image than a high-capacity MLP. The feature extractor outputs a geometry representation (*i.e.* the feature map) of the full scene in a single forward pass. During training, it receives gradient information from all points sampled in the camera frustum, conditioned on the input image. Thus, potential noise from small visual details gets averaged out. On the other hand, the MLP out-

puts density based on the coordinates and is conditioned on a local feature. The coordinates and feature are different for every sampled point, rather than per scene. Consequently, noise will affect the MLP training significantly more.

Introducing the sampling of color from the input frames further boosts accuracy, especially for occupancy estimation in occluded areas. We hypothesize that only predicting density simplifies the training task significantly. Crucially, the network does not have to hallucinate colors in occluded regions. Additionally, color sampling enforces strict multi-view consistency. The network cannot compensate for imperfect geometry by predicting the correct color.

Finally, the results show that our policy of discarding invalid rays during training improves accuracy by reducing noise in the training signal. This mainly affects the border regions of the frustum.

## 5. Conclusion

In this paper, we introduced a new approach to learning to estimate the 3D geometric structure of a scene from a single image. Our method predicts a continuous density field, which can be evaluated at any point in the camera frustum. The key contributions in our paper are 1. color sampling, 2. architecture improvements, and 3. a new self-supervised loss formulation. This enables us to train a network on large in-the-wild datasets with challenging scenes, such as KITTI, KITTI-360, and RealEstate10K. We show that our method is able to capture geometry in occluded areas. We evaluate depth maps synthesized from the predicted representation achieving comparable results to state-of-the-art methods. Despite only predicting geometry, our model even achieves high accuracy for novel view synthesis from a single image. Finally, we justify all of our design choices through detailed ablation studies.

**Acknowledgements.** This work was supported by the ERC Advanced Grant SIMULACRON, by the Munich Center for Machine Learning and by the EPSRC Programme Grant VisualAI EP/T028572/1. C. R. is supported by VisualAI EP/T028572/1 and ERC-UNION-CoG-101001212.

## References

- [1] Shubhra Aich, Jean Marie Uwabeza Vianney, Md Amirul Islam, and Mannat Kaur Bingbing Liu. Bidirectional attention network for monocular depth estimation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11746–11752. IEEE, 2021. 2
- [2] 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. 2
- [3] 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. 2
- [4] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multi-map slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021. 2
- [5] Anh-Quan Cao and Raoul de Charette. Scenerf: Self-supervised monocular 3d scene reconstruction with radiance fields. *arXiv*, 2022. 3
- [6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *arXiv preprint arXiv:2203.09517*, 2022. 3
- [7] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022. 2
- [8] Helisa Dharmo, Nassir Navab, and Federico Tombari. Object-driven multi-layer scene decomposition from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5369–5378, 2019. 3
- [9] Helisa Dharmo, Keisuke Tateno, Iro Laina, Nassir Navab, and Federico Tombari. Peeking behind objects: Layered depth prediction from a single image. *Pattern Recognition Letters*, 125:333–340, 2019. 3
- [10] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 2, 7
- [11] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2002–2011, 2018. 2
- [12] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 2, 5
- [13] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 270–279, 2017. 2
- [14] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3828–3838, 2019. 1, 2, 4, 5, 6, 7
- [15] Juan Luis GonzalezBello and Munchurl Kim. Forget about the lidar: Self-supervised depth estimators with med probability volumes. *Advances in Neural Information Processing Systems*, 33:12626–12637, 2020. 1, 2
- [16] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raventos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2485–2494, 2020. 1, 2, 7
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [18] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. 2
- [19] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 3
- [20] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022. 2
- [21] Yevhen Kuznietsov, Jorg Stuckler, and Bastian Leibe. Semi-supervised deep learning for monocular depth map prediction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6647–6655, 2017. 2
- [22] Sihaeng Lee, Janghyeon Lee, Byungju Kim, Eojindil Yi, and Junmo Kim. Patch-wise attention network for monocular depth estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1873–1881, 2021. 2
- [23] Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu Wang, and Gim Hee Lee. Mine: Towards continuous depth mpi with nerf for novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12578–12588, 2021. 3, 6, 7, 8
- [24] Zhenyu Li, Zehui Chen, Xianming Liu, and Junjun Jiang. Depthformer: Exploiting long-range correlation and local information for accurate monocular depth estimation. *arXiv preprint arXiv:2203.14211*, 2022. 2
- [25] Zhenyu Li, Xuyang Wang, Xianming Liu, and Junjun Jiang. Binsformer: Revisiting adaptive bins for monocular depth estimation. *arXiv preprint arXiv:2204.00987*, 2022. 2
- [26] Yiyi Liao, Jun Xie, and Andreas Geiger. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2, 5
- [27] Fayao Liu, Chunhua Shen, Guosheng Lin, and Ian Reid. Learning depth from single monocular images using deep

- convolutional neural fields. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):2024–2039, 2015. [2](#)
- [28] 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. [3](#)
- [29] Chenxu Luo, Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, Ram Nevatia, and Alan Yuille. Every pixel counts++: Joint learning of geometry and motion with 3d holistic understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2624–2641, 2019. [1](#), [2](#), [7](#)
- [30] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):99–108, 1995. [2](#), [3](#)
- [31] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. [3](#)
- [32] 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](#)
- [33] Norman Müller, Andrea Simonelli, Lorenzo Porzi, Samuel Rota Bulò, Matthias Nießner, and Peter Kotschieder. Autorf: Learning 3d object radiance fields from single view observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3971–3980, 2022. [3](#)
- [34] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. [3](#)
- [35] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. [2](#)
- [36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. [5](#)
- [37] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision*, pages 523–540. Springer, 2020. [3](#)
- [38] Barbara Roessle, Jonathan T Barron, Ben Mildenhall, Pratul P Srinivasan, and Matthias Nießner. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12892–12901, 2022. [2](#)
- [39] Mehdi SM Sajjadi, Henning Meyer, Etienne Pot, Urs Bergmann, Klaus Greff, Noha Radwan, Suhani Vora, Mario Lučić, Daniel Duckworth, Alexey Dosovitskiy, et al. Scene representation transformer: Geometry-free novel view synthesis through set-latent scene representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6229–6238, 2022. [3](#)
- [40] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [2](#)
- [41] 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](#)
- [42] Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski. Layered depth images. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pages 231–242, 1998. [3](#)
- [43] Prafull Sharma, Ayush Tewari, Yilun Du, Sergey Zakharov, Rares Ambrus, Adrien Gaidon, William T Freeman, Fredo Durand, Joshua B Tenenbaum, and Vincent Sitzmann. Seeing 3d objects in a single image via self-supervised static-dynamic disentanglement. *arXiv preprint arXiv:2207.11232*, 2022. [2](#), [3](#)
- [44] Chang Shu, Kun Yu, Zhixiang Duan, and Kuiyuan Yang. Feature-metric loss for self-supervised learning of depth and egomotion. In *European Conference on Computer Vision*, pages 572–588. Springer, 2020. [1](#), [2](#), [7](#)
- [45] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European conference on computer vision*, pages 746–760. Springer, 2012. [2](#)
- [46] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 175–184, 2019. [3](#)
- [47] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11358–11367, 2021. [3](#)
- [48] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 551–560, 2020. [3](#), [8](#)
- [49] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3d scene inference via view synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 302–317, 2018. [3](#), [7](#), [8](#)
- [50] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [4](#)

- [51] Jamie Watson, Michael Firman, Gabriel J Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2162–2171, 2019. [1](#), [2](#), [7](#)
- [52] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7467–7477, 2020. [3](#), [8](#)
- [53] Felix Wimbauer, Shangzhe Wu, and Christian Rupprecht. De-rendering 3d objects in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18490–18499, 2022. [3](#)
- [54] Felix Wimbauer, Nan Yang, Lukas Von Stumberg, Niclas Zeller, and Daniel Cremers. Monorec: Semi-supervised dense reconstruction in dynamic environments from a single moving camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6112–6122, 2021. [2](#)
- [55] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1–10, 2020. [3](#)
- [56] Nan Yang, Rui Wang, Jorg Stuckler, and Daniel Cremers. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 817–833, 2018. [2](#)
- [57] 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](#), [3](#), [6](#), [7](#), [8](#)
- [58] Weihao Yuan, Xiaodong Gu, ZuoZhuo Dai, Siyu Zhu, and Ping Tan. New crfs: Neural window fully-connected crfs for monocular depth estimation. *arXiv preprint arXiv:2203.01502*, 2022. [1](#), [2](#)
- [59] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 340–349, 2018. [1](#), [2](#)
- [60] Kaichen Zhou, Lanqing Hong, Changhao Chen, Hang Xu, Chaoqiang Ye, Qingyong Hu, and Zhenguo Li. Devnet: Self-supervised monocular depth learning via density volume construction. *arXiv preprint arXiv:2209.06351*, 2022. [2](#), [7](#)
- [61] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1851–1858, 2017. [1](#), [2](#)
- [62] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. [3](#), [5](#), [7](#)
- [63] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *European conference on computer vision*, pages 286–301. Springer, 2016. [3](#)