# Depth-supervised NeRF: Fewer Views and Faster Training for Free

Kangle Deng[1]       Andrew Liu[2]       Jun-Yan Zhu[1]       Deva Ramanan[1,3]
[1]Carnegie Mellon University       [2]Google       [3]Argo AI
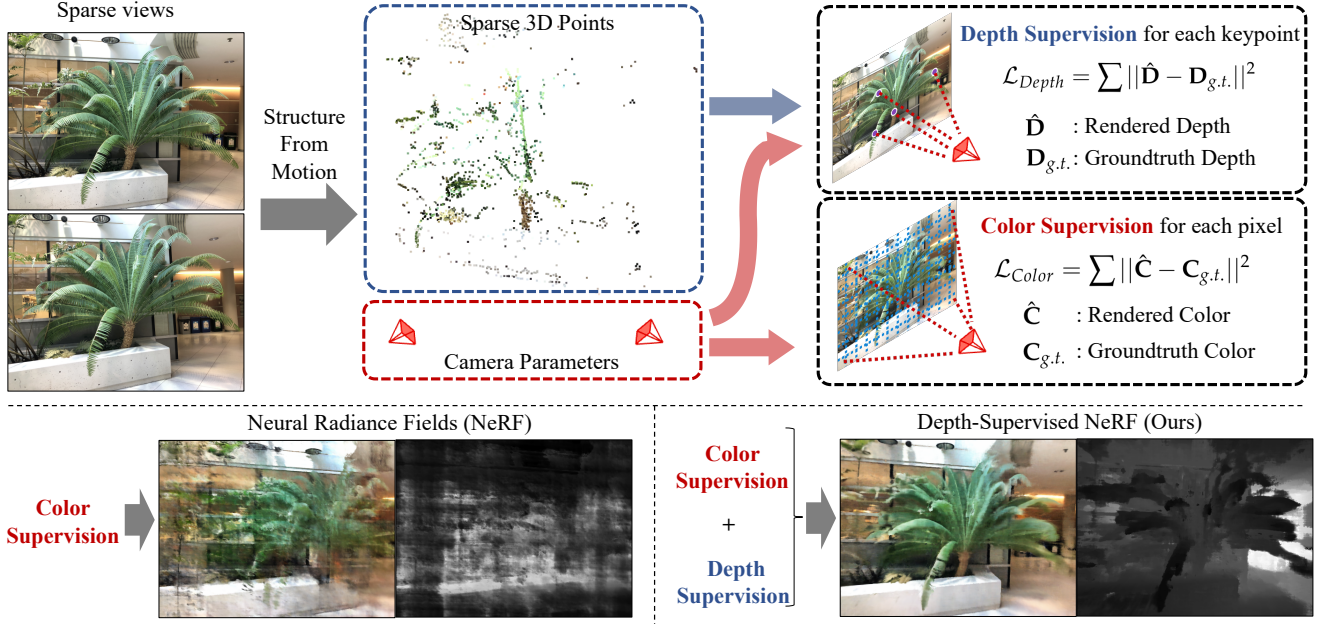
Figure 1: Because NeRFs can be difficult to train given insufficient color supervision (i.e., too few input images), we make use of additional supervision from depth. One such source is 3D point clouds computed when running structure-from-motion (SFM) to estimate camera poses. We point out that such readily-available sparse 3D signals can be used as additional free supervision. We impose a loss to ensure that depth rendered along rays that intersect these 3D points are close to their observed depth. Because depth supervision is complementary to many learning-based NeRF pipelines, it can be combined with such approaches to dramatically reduce overfitting and speed up training.

## Abstract

*One common failure mode of Neural Radiance Field (NeRF) models is fitting incorrect geometries when given an insufficient number of input views. We propose DS-NeRF (Depth-supervised Neural Radiance Fields), a loss for learning neural radiance fields that takes advantage of readily-available depth supervision. Our key insight is that sparse depth supervision can be used to regularize the learned geometry, a crucial component for effectively rendering novel views using NeRF. We exploit the fact that current NeRF pipelines require images with known camera poses that are typically estimated by running structure-from-motion (SFM). Crucially, SFM also produces sparse 3D points that can be used as "free" depth supervision during training: we simply add a loss to ensure that depth rendered along rays that intersect these 3D points is close to the observed depth. We find that DS-NeRF can render more accurate images given fewer training views while training 2-6x faster. With only two training views on real-world images, DS-NeRF significantly outperforms NeRF as well as other sparse-view variants. We show that our loss is compatible with these NeRF models, demonstrating that depth is a cheap and easily digestible supervisory signal. Finally, we show that DS-NeRF supports other types of depth supervision such as scanned depth sensors and RGBD reconstruction outputs.*

## 1. Introduction

Neural rendering with implicit representations has become a widely-used technique for solving many kinds of vision and graphics tasks ranging from view synthe-
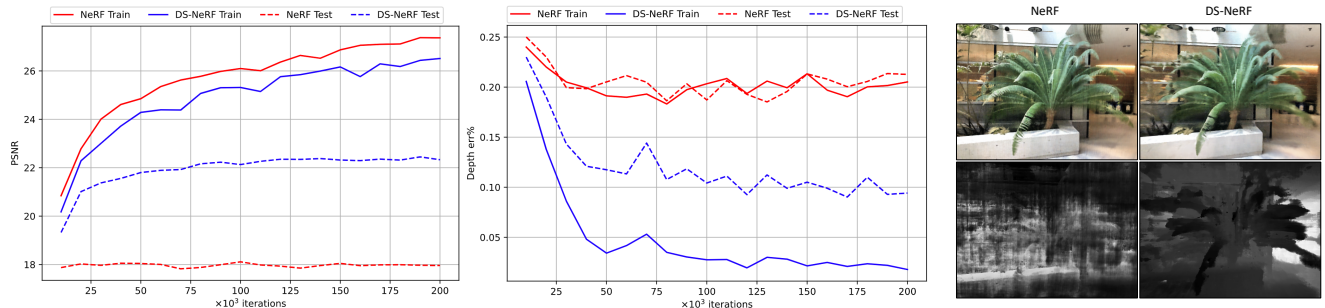
Figure 2: **NeRF can overfit** to a small number of training views (here, two), as shown by the PSNR gap between train and test view renderings (left). This overfitting is due to wildy inccurate geometries being learned, as evidenced by the large error of the rendered depth maps, even for training views (middle). We visualize both the rendered image and the rendered depth map for a training view on the right. DS-NeRF exploits (sparse) depth supervision to learn far more accurate scene geometry, resulting in better generalization to novel views.

sis [28, 18], to re-lighting [16, 14], to pose and shape estimation [19, 24, 36], to 3D-aware image synthesis and editing [26, 1, 11], to modeling dynamic scene [20, 22, 9]. These prior works have demonstrated the powerful, yet extraordinarily flexible, ability of implicit functions to model complex scenes and 3D concepts.

Among them is the seminal work of Neural Radiance Fields (NeRF) [18]. Mildenhall *et al*. demonstrated impressive view synthesis results by using implicit functions to encode volumetric density and color observations before rendering arbitrary views with standard ray-tracing. Since then, the NeRF representation has been extended to learning more complex phenomenons like dynamic scene flows [35, 20] and illumination effects [29, 14].

In spite of this, NeRF has several limitations. Reconstructing both the scene appearance and geometry can be ill-posed given a small number of input views. Figure 2 shows that NeRF can learn wildly inaccurate scene geometries that still accurately render train-views. However, such models produce poor renderings of novel test-views, essentially overfitting to the train set. Moreoever, even given a large number of input views, NeRF can still be time-consuming to train; it often takes between ten hours to several days to model a single scene at moderate resolutions on a single GPU. The training is slow due to both the expensive ray-casting operations and lengthy optimization process.

In this work, we explore depth as an additional, cheap source of supervision to regularize the geometry learned by NERF and improve NeRF's training. To train a NeRF, the camera parameters of each view are often first estimated to determine how rays will be casted through a volumetric representation. This is commonly solved using structure-from-motion (SFM) solvers like COLMAP [25]. In addition to camera pose and intrinsics, COLMAP also outputs two additional signals, sparse 3D point clouds and the reprojection errors between detected 2D keypoints and projected 3D points; *surprisingly these outputs are ignored by NeRF*.

We posit that additional depth supervision provides valuable training gradients that are currently missing from pixel reconstruction losses.

Our work makes use of this "free" supervision through the use of a depth supervision loss. In addition to integrating the color along a ray, one can also integrate the occupancy to compute an expected depth from a particular camera viewpoint. By sampling rays that intersect with the sparse 3D keypoints detected by SFM, we can impose a loss between the SFM-estimated depth and the one predicted by NeRF. This is a significantly stronger signal than reconstructing RGB pixels. Without depth supervision, NeRF is implicitly solving a 3D correspondence problem between multiple views. However, the sparse version of this exact problem has already been solved by structure-from-motion, whose solution is given by the sparse 3D keypoints. Therefore depth supervision improves NeRF by anchoring its search over implicit correspondences with sparse explicit ones.

Further, we find that DS-NeRF can utilize other sources of depth supervision besides structure-from-motion like depth sensors, which are becoming more ubiquitous. Because our approach produces high quality results with comparatively lighter compute constraints, DS-NeRF opens up potential NeRF-based applications with small form-factor camera arrays (e.g., multi-camera phones) or edge compute sensors (e.g., traffic sensors).

Our experiments show that this surprisingly simple idea translates to massive improvements in training NeRFs and its variations, regarding both the training speed and the amount of training data needed. We observe that depth-supervised NeRF can accelerate model training by 2-6x while producing results with the same quality. For sparse view settings, experiments show that our method synthesizes better results compared to the original NeRF and recent sparse-views NeRF models [37, 30] on both NeRF Real-world [18] and Redwood-3dscan [3] We show that our depth supervision loss works well with depth derived from other sources such

as a depth camera and RGBD reconstruction outputs. As a result, our loss can be seamlessly integrated into existing NeRF models. Our code and more results are available at our GitHub repo and website.

## 2. Related Work

**Implicit Representations for 3D.** Devising efficient 3D representations has been a long-standing research problem. Recent methods have used neural networks to encode and decode 3D information via implicit functions [19, 2, 15, 24]. These methods take a set of 2D images or 3D shapes to learn an implicit function that maps coordinates to surface and texture properties at those coordinates. Later works explored learning implicit functions without relying on explicit 3D supervision [28, 18].

NeRFs [18] learn an implicit function to model a volumetric radiance field using standard volumetric rendering [8] and alpha compositing techniques [21]. While NeRF is capable of rendering new views with impressive visual qualities, one drawback is the large number of training views and lengthy optimization required to learn the underlying geometry correctly. Followup works have explored ways to improve NeRF training by reducing the number of views required [37] or increasing the speed of convergence [30]. Unlike these works which derive performance gains by learning category priors, our approach uses category-free supervision in the form of readily available 3D observation to directly supervise the volumetric geometry, resulting in a more data-efficient and faster NeRF training process.

**Novel View Synthesis.** Prior to implicit representations, popular approaches for view synthesis used explicit representations like multiplane images [40, 5, 17, 31], meshes [27, 6], or point clouds [34] to render novel views. While such approaches are able to render quickly at test time, they have limited expressiveness with their representation.

NeRF uses a test-time optimization procedure to generate novel views of single scenes and is capable of fitting to image sets without relying on category-specific or domain priors. Because there are no priors, the optimization for learning NeRF is slow and data-inefficient. Tancik et al. [30] propose speeding up the training behavior of NeRF by using meta-learning [4] on a domain of training scenes to find a good implicit function initialization. Another approach by Yu et al. [37] trains an encoder on scenes such that test-time optimization is no longer needed. However, these approaches are limited in that they are sensitive to priors in the training data of a certain category (e.g., chairs) that may not translate well to unseen domains (e.g., real-world images). In contrast, 2D keypoints and 3D correspondence are geometric properties of any 3D scene. Therefore, regardless of the input scene and its domain, we can always use depth supervision to achieve better and faster training.

**Depth and Beyond.** Depth is a valuable source of data for view synthesis and 3D representations. Early works have focused on recovering depth via algorithms like multiview stereo and SFM [25]. While these tools have become staples of 3D learning, recent works have focused on using deep learning to extract dense depth in the monocular setting [10, 23] or from videos [12]. These works have made depth readily accessible for recent 3D photography method [27] and video-to-video synthesis system [13].

In our work, we propose to incorporate sparse depth information derived from SFM to assist in learning NeRF-based view synthesis. While such a depth loss has been previously explored in an MPI setting [31], we find its addition to NeRF significantly improves training and enables previously impossible NeRF applications.

## 3. Depth-supervised NeRF (DS-NeRF)

We now present our Depth-supervised Neural Radiance Fields. We first go over the volumetric rendering used by NeRF-based methods. Then we introduce sampling rays with depth information. Finally, we present our depth-supervision loss over these sampled rays.

### 3.1. Volumetric Rendering

A Neural Radiance Field takes a set of images and encodes a scene as a volume density and emitted radiance. More specifically, for a given 3D point $\mathbf{x} \in \mathbb{R}^3$ and a particular viewing direction $\mathbf{d} \in \mathbb{R}^3$, NeRF learns an implicit function $f$ that estimates the differential density $\sigma$ and RGB color $\mathbf{c}$ like so: $f(\mathbf{x}, \mathbf{d}) = (\sigma, \mathbf{c})$.

To render a 2D image given a pose $\mathbf{P}$, we shoot rays $\mathbf{r}$ originating from the $\mathbf{P}$'s center of projection $\mathbf{o}$ in direction $\mathbf{d}$ derived from its intrinsics. We sample the implicit radiance field at coordinates lying along this ray to approximate its color:

$$\hat{\mathbf{C}} = \int_{t_n}^{t_f} T(t)\sigma(t)\mathbf{c}(t)dt, \tag{1}$$

where $T(t) = \exp(-\int_{t_n}^{t} \sigma(s)ds)$ checks for occlusions by integrating the differential density over points between $t_n$ to $t$. $t$ parameterizes the aforementioned ray as $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$. We evaluate rays on this implicit function between near $(t_n)$ and far $(t_f)$ bounds.

The final NeRF rendering loss is given by a reconstruction loss over colors returned from rendering the set of rays $\mathcal{R}(\mathbf{P})$ produced by a particular camera parameter $\mathbf{P}$.

$$\mathcal{L}_{\text{Color}} = \sum_{\mathbf{r} \in \mathcal{R}(\mathbf{P})} \left\| \hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2. \tag{2}$$

One limitation of NeRF is its inefficient training, both in terms of the number of test-time optimizations required to fit a scene as well as the number of input images. Because supervision is only derived from the image data itself, learning

a volumetric geometry requires shooting rays and sampling points to trivially learn that the differential density for empty regions is zero.

These problems are present in any kind of ray-based techniques like Marching Cubes [28] and NeRF's ray-tracing. While certain optimizations like hierarchical sampling improve sample efficiency, test-time optimization over few views will still be problematic, especially when compared to prior view synthesis methods that are capable of inferring novel views from a single image [34, 31, 27].

We posit that while NeRF's are effective at learning arbitrarily complex radiance fields via their volumetric rendering, many scenes can be represented using a simpler geometric representation (point clouds, meshes, multi-plane images, depth maps, etc). Our goal is to unify NeRF's powerful implicit rendering capabilities with the computing efficiency of explicit 3D reconstruction methods through the use of explicit geometric supervision.

### 3.2. Motivating Explicit Supervision

Let's consider the training of a newly initialized NeRF. Without any priors, the differential density is uniformly allocated across the volume. For learning arbitrary radiance fields, a uniform volume is a good choice because it's the median initialization. However this is also a poor initialization as in typical scenes, the majority of space is either occupied or empty.

A poor initialization can result in longer test-time optimizations and degenerate behaviors. For example given only two views, NeRF could achieve perfect rendering by creating a volumetric wall at the cameras' near bounds and returning the image's pixels as the color.

One way to resolve this ambiguity is to use more views. Additional views provide sufficient supervision to prevent the network from collapsing. However in many cases, we do not have extra views. Therefore we consider alternative sources of supervision like **depth** which can force the network to explain a volumetric representation that is consistent with both pixels as well as measurements derived from structure-from-motion.

### 3.3. Deriving Depth

Collections of images rarely come annotated with ground truth depth or camera poses, especially for imagery collected in an uncontrolled environment. In order to train a NeRF over these images, camera poses are estimated using structure-from-motion (SFM) frameworks like COLMAP [25] to simultaneously estimate camera pose $\mathbf{P}$ and intrinsics $\mathbf{K}$ from a collection of imagery.

In addition to estimating camera parameters $(\mathbf{P_1 K_1}, \mathbf{P_2 K_2}, \ldots)$, structure-from-motion also estimates 3D keypoints $\{\mathbf{X} : \mathbf{x}_1, \mathbf{x}_2, \ldots \in \mathbb{R}^3\}$ across multiple views. Here $\mathbf{x}_i$ is a detected keypoint in a global 3D

coordinate space. Importantly, SFM also models scene occlusion via bundle adjustment, explicitly identifying subsets set of keypoints visible from a particular camera $j$: $\mathbf{X}_j \subset \mathbf{X}$.

### 3.4. Depth Supervision Loss

Depth images from a particular pose can also be rendered in a similar fashion to rendering RGB pixels. We render depth using the initial formulation proposed in the original NeRF paper. Given a ray that is parameterized as $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

$$\hat{\mathbf{D}}(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(t)t dt, \tag{3}$$

where $\hat{\mathbf{D}}(\mathbf{r})$ is the expected depth along the camera axis of ray $\mathbf{r}$ and $T(t)$ is the same occlusion checking integration defined earlier when rendering color in Eq 1.

We propose using the 3D keypoints provided by SFM to supervise the depth of specific rays. We start with the NeRF training formulation, given an image $I_j$ and estimated camera parameters $\mathbf{P_j}$, $\mathbf{K_j}$, we would like to render a set of its camera's rays. Unlike RGB supervision which is dense, SFM only produces sparse depth supervision in terms of visible keypoints $\mathbf{x}_i \in \mathbf{X}_j$. For each keypoint $i$, we can cast a ray to camera $j$, written as $\mathbf{r}_{ij}(t) = \mathbf{o}_j + \mathcal{F}_j(\mathbf{x}_i - \mathbf{o}_j)t$ which samples points lying between camera $j$'s center of projection $\mathbf{o}_j$ and keypoint $\mathbf{x}_i$. $\mathcal{F}_j$ is a function which scales the ray to have length 1 along $j$'s camera axis. In practice we implement this by using the NeRF ray corresponding to $\mathbf{x}_i$'s 2D image coordinates in image $j$.

To calculate depth of keypoints in camera $j$, we reproject them using $\mathbf{P}_j$ before projecting the result onto its unit camera axis $[0, 0, 1]$. When applicable, we implicitly convert keypoints to their homogenous representation. Our final depth-supervised loss is:

$$\mathcal{L}_{\text{Depth}} = \sum_{\mathbf{x}_i \in \mathbf{X}_j} w_i \left| \hat{\mathbf{D}}(\mathbf{r}_{ij}) - (\mathbf{P}_j \mathbf{x}_i) \cdot [0, 0, 1] \right|^2, \tag{4}$$

where $w_j$, described below, adaptively weighs keypoints.

**Confidence-weighed Keypoints.** SFM might produce spurious correspondences and/or unreliable keypoints that could hinder depth supervision. To reduce the influence of such unreliable keypoints, we weight each keypoint by its reprojection error estimated by SFM. Specifically, given a 3D keypoint $\mathbf{x}_i$ visible in camera $j$, the reprojection error $e_{ij}$ is the distance in pixels between projected image coordinates $\mathbf{K}_j \mathbf{P}_j \mathbf{x}_i$ and the detected keypoint in 2D. We use $\mathbf{x}_i$'s total reprojection error of $e_i = \sum_j e_{ij}$ to measure the confidence weight $w_i$ of a particular keypoint:

$$w_i = \exp\left(-(\frac{e_i}{\bar{e}})^2\right), \tag{5}$$

where $\bar{e}$ is the average absolute error over all keypoints in a scene.

**Training loss.** Our total training loss is:

$$\mathcal{L} = \mathcal{L}_{\text{Color}} + \lambda_D \mathcal{L}_{\text{Depth}}, \tag{6}$$

where $\lambda_D$ is a hyper-parameter balancing color and depth supervision.

**Free training.** So far we have discussed ways to extract free supervision, now we show how this additional supervision can be seamlessly integrated into training with minimal extra compute. Although depth rendering and color rendering have different integral formulas, both require ray-tracing $\mathbf{r}(t)$ as input. By recognizing that RGB colors are a dense signal, we can sample off-center pixel rays and compute their color supervision via bilinear interpolation of the input image. This allows us to turn depth supervised rays into depth-and-color supervised rays, simultaneously training RGB and depth in a single forward pass.

## 4. Experiments

We first evaluate the input data efficiency on view synthesis over several datasets in Section 4.3. For relevant NeRF-related methods, we also evaluate the error of rendered depth maps in Section 4.4. Finally, we analyze training speed improvements in Section 4.5.

### 4.1. Datasets

We use several multi-view datasets to evaluate different approaches on various tasks.

**NeRF Real-world Data (NeRF Real)** [17, 18] contains eight complex, real world scenes captured with a forward-facing camera. To evaluate view synthesis under different view sparsity, we create subsets of training images for each scene of sizes 2, 5, and 10 views. For every subset, we run COLMAP [25] over its training images to estimate cameras and collect sparse keypoints for depth supervision.

**DTU MVS Dataset (DTU)** [7] captures various objects in a controlled environment from multiple perspectives. Following Yu *et al*.'s setup in PixelNeRF [37], we evaluated on the same 15 test scenes and views. For each scene, we used their subsets of size 3, 6, 9 training views. To get sparse point clouds, we run COLMAP combined with the ground truth calibrated camera poses. Images are down-sampled to a resolution of $400 \times 300$ for training and evaluation.

**Redwood-3dscan (Redwood)** [3] contains RGB-D videos of various categories such as cars, chairs, and plants. We select 5 RGB-D sequences and create subsets of 2, 5, and 10 training frames for each object. We run COLMAP to get their camera poses and sparse point clouds. In addition, Redwood provides a depth map collected from a depth sensor. This allows us to compute each NeRF model's depth error by comparing against a ground truth depth measurement. In order to connect the scale of COLMAP's pose with the

scanned depth, we apply an affine transformation so the scanned depth best fits detected keypoints by solving a least-square optimization. Please see our appendix for additional details when processing the scanned depth.

### 4.2. Comparisons

First we compare with non-NeRF based approaches like Local Lightfield Fusion (*LLFF*) [17], an MPI-based representation that fuses learned representations from multiple view points. Next we consider a set of NeRF baselines. More details can be found in Appendix A.3.

**PixelNeRF** [37] expands upon NeRF by using an encoder to train a general model across multiple scenes. For pixelNeRF, we consider 2 variations. *pixelNeRF-DTU* is evaluated using their pre-trained checkpoint on DTU. For cases where the train and test domain are different, we fine-tune additional iterations on each test scene to get *pixelNeRF finetuned*.

**MetaNeRF** [30] finds a better NeRF initialization over a domain of training scenes before running test-time optimization on new scenes. Because DTU is the only dataset large enough for meta-learning, we only consider the *metaNeRF-DTU* baseline which learns an initialization over DTU for 40K meta-iterations (equivalent to 40K $\times$ 64 inner NeRF optimizations) and then finetunes for 1000 steps on new scenes. We purposefully follow metaNeRF's ShapeNet setup to demonstrate its susceptibility to differences between training and testing domains, however we explore finetuning metaNeRF for longer in the supplementary.

**IBRNet** [32] extends NeRF by using a MLP and ray transformer to estimate radiance and volume density.

As our depth supervision loss does not require additional annotation or data collection, our loss readily fits in a wide range of NeRF-based architectures. Here, we also incorporate our loss into pixelNeRF and IBRNet with little effort. Finally, we consider an ablation of our full model *w/o confidence-weighted loss*.

### 4.3. Few-input view synthesis

We start by comparing each method on rendering test views from few inputs. For view synthesis, we report three metrics (PSNR, SSIM [33], and LPIPS [39]) that evaluate the quality of rendered views against a ground truth.

**NeRF Real.** As seen in Table 1 on the NeRF Real, our approach produces better metrics than NeRF and LLFF, especially when only two and fives input views are available. We present qualitative comparisons in Figure 3. metaNeRF-DTU and pixelNeRF-DTU visually struggle because DTU contains priors that cannot be transferred here. While finetuning of pixelNeRF-DTU helps, finetuning with depth-supervision can further improve results.

**DTU.** We also show numerical performance on DTU in Table 3 and qualitative results in Figure 4. We find that DS-

5

| NeRF Real | PSNR↑ | | | SSIM↑ | | | LPIPS↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view |
| LLFF | 16.3 | 18.4 | 22.3 | 0.48 | 0.49 | 0.53 | 0.55 | 0.51 | 0.53 |
| NeRF | 18.2 | 22.6 | 24.2 | 0.50 | 0.57 | 0.58 | 0.52 | 0.50 | 0.52 |
| metaNeRF-DTU | 13.1 | 13.8 | 14.3 | 0.43 | 0.45 | 0.46 | 0.89 | 0.88 | 0.87 |
| pixelNeRF-DTU | 12.7 | 12.4 | 12.8 | 0.39 | 0.40 | 0.40 | 0.82 | 0.87 | 0.81 |
|    finetuned | 20.2 | 23.1 | 24.1 | 0.56 | 0.59 | 0.63 | 0.53 | 0.53 | 0.41 |
|    finetuned with DS | 20.9 | 23.1 | 24.4 | 0.54 | 0.61 | 0.66 | 0.55 | 0.47 | 0.42 |
| IBRNet | 18.4 | 22.9 | 24.3 | 0.50 | 0.51 | 0.54 | 0.53 | 0.54 | 0.51 |
|    finetuned with DS | 22.2 | 24.3 | **25.2** | 0.63 | 0.66 | 0.68 | 0.39 | **0.36** | 0.38 |
| DS-NeRF | **22.3** | **24.5** | **25.2** | **0.68** | **0.69** | 0.69 | **0.38** | 0.37 | 0.37 |
|    w/o weighted loss | 21.9 | 24.2 | 25.1 | 0.65 | **0.69** | **0.71** | 0.39 | **0.36** | **0.36** |

Table 1: **View Synthesis on NeRF Real:** We evaluate view synthesis quality for various methods when given 2, 5, 10 views from NeRF Real. Compared with both NeRF and non-NeRF based methods, DS-NeRF renders higher quality test views. On the other hand, metaNeRF-DTU [30] and pixelNeRF-DTU [37], methods that were trained with the DTU dataset [7], produce poor quality results due to the domain gap. Additionally we apply additional finetuning of pixelNeRF-DTU on scenes from NeRF Real and find that employing depth supervision can improve results. Likewise we observe similar trends when adding depth supervision to finetuning IBRNet [32].

| Depth err%↓ | NeRF real-world | | | Redwood-3dscan | | |
|---|---|---|---|---|---|---|
| | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view |
| NeRF | 20.32 | 15.00 | 12.41 | 25.32 | 24.34 | 21.34 |
| metaNeRF-DTU | 22.23 | 22.07 | 22.30 | 20.84 | 21.12 | 20.96 |
| pixelNeRF-DTU | 22.12 | 22.09 | 22.06 | 19.46 | 19.87 | 19.54 |
| DS-NeRF | **10.41** | **8.61** | **8.15** | 11.42 | 10.43 | 9.43 |
| DS-NeRF w/ RGB-D | - | - | - | **5.81** | **5.31** | **4.22** |

Table 2: **Depth Error:** We compare rendered depth to the "ground-truth" obtained from either COLMAP[25] (run on *all* available views in NeRF real-world) or Redwood[3] RGB-D reconstruction obtained on depth scans from Redwood. DS-NeRF produces volumes that better represent the underlying geometry as indicated by the lower depth errors of test views. Training with the RGB-D reconstruction output further improves performance. We report the average error between up-to-scale rendered and reference depths, normalized by the reference.

| DTU | PSNR↑ | | | SSIM↑ | | | LPIPS↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 3-view | 6-view | 9-view | 3-view | 6-view | 9-view | 3-view | 6-view | 9-view |
| NeRF | 9.85 | 18.59 | 22.14 | 0.37 | 0.72 | **0.82** | 0.62 | 0.35 | 0.26 |
| metaNeRF-DTU | 18.19 | 18.75 | 20.15 | 0.60 | 0.61 | 0.67 | 0.40 | 0.41 | 0.35 |
| pixelNeRF-DTU | **19.33** | 20.43 | 21.10 | **0.70** | **0.73** | 0.76 | **0.39** | 0.36 | 0.34 |
| DS-NeRF | 16.63 | 20.52 | **22.26** | 0.56 | 0.72 | 0.79 | 0.45 | 0.33 | **0.24** |
|    w/o weighted loss | 16.52 | **20.54** | 22.23 | 0.54 | **0.73** | 0.77 | 0.48 | **0.31** | 0.26 |

Table 3: **View Synthesis on DTU:** We evaluate on 3, 6, and 9 views respectively for 15 test scenes from the DTU dataset [7]. pixelNeRF-DTU [37] and metaNeRF-DTU [30] perform well given that the domain overlap between training and testing, although DS-NeRF is still competitive on view synthesis for 6 and 9 views. We explore variants of pixelNeRF augmented with depth supervision in Table 4.
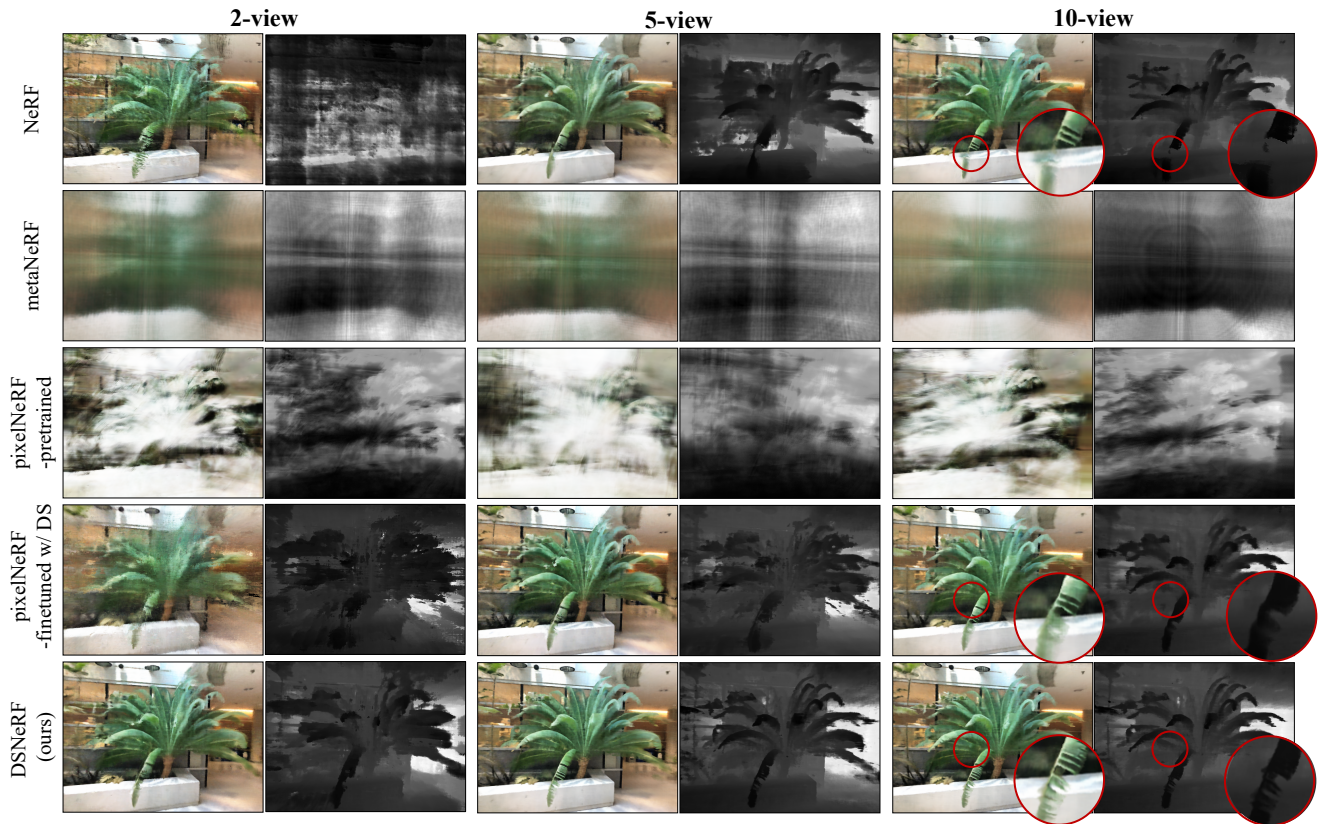
Figure 3: **Qualitative Comparison on NeRF Real:** We render image and depth from a test view for various NeRF models trained on 2, 5, and 10 views. We find that methods trained with DTU struggle on NeRF Real. On the other hand, methods using depth-supervision are able to render test views with realistic geometry, even when only 2 views are provided. Please refer to Table 1 for quantitative comparison.

.

| DTU (PixelNeRF w/ DS) | PSNR↑ | | | SSIM↑ | | | LPIPS↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 3-view | 6-view | 9-view | 3-view | 6-view | 9-view | 3-view | 6-view | 9-view |
| pixelNeRF-DTU | 19.33 | 20.43 | 21.10 | 0.70 | **0.73** | 0.76 | 0.39 | 0.36 | 0.34 |
|     trained with DS | **19.87** | **20.64** | **21.25** | **0.71** | **0.73** | **0.77** | **0.37** | **0.35** | **0.32** |
| pixelNeRF-DTU finetuned | 19.97 | 20.89 | 22.32 | **0.72** | 0.74 | **0.79** | **0.37** | 0.34 | 0.33 |
|     finetuned with DS | **20.01** | **21.07** | **22.42** | **0.72** | **0.76** | **0.79** | **0.37** | **0.32** | **0.32** |

Table 4: **PixelNeRF[37] and Depth Supervision:** Depth supervision can be combined with recent NeRF-based approaches with little effort. We incorporate depth supervision during either pixelNeRF training or finetuning on test scenes. Adding depth supervision during either stage consistently improves performance over the original version.

NeRF renders images from 6 and 9 input views that are competitive with pixelNeRF-DTU, however metaNeRF-DTU and pixelNeRF-DTU are able to outperform DS-NeRF on 3-views. This is not particularly surprising as both methods are trained on DTU scenes and therefore can fully leverage dataset priors. However this also highlights their apparent weakness. As evident by their performances on NeRF Real and Redwood, these DTU-pretrained models struggle to perform well outside of DTU. Our full approach is capable of achieving good rendering results regardless of domain.

**Redwood.** Like NeRF Real, we find similar improvements in performance across the Redwood dataset in Table 5. Because Redwood includes depth measurements collected with a sensor, we also consider how alternative sources of depth supervision can improve results. We train a DS-NeRF, replacing COLMAP supervision with the scaled Redwood depth measurements. We find that the denser depth helps DS-NeRF even more, achieving a PSNR scoreof 22.42 on 2-views which is competitive with NeRF on 5-views.

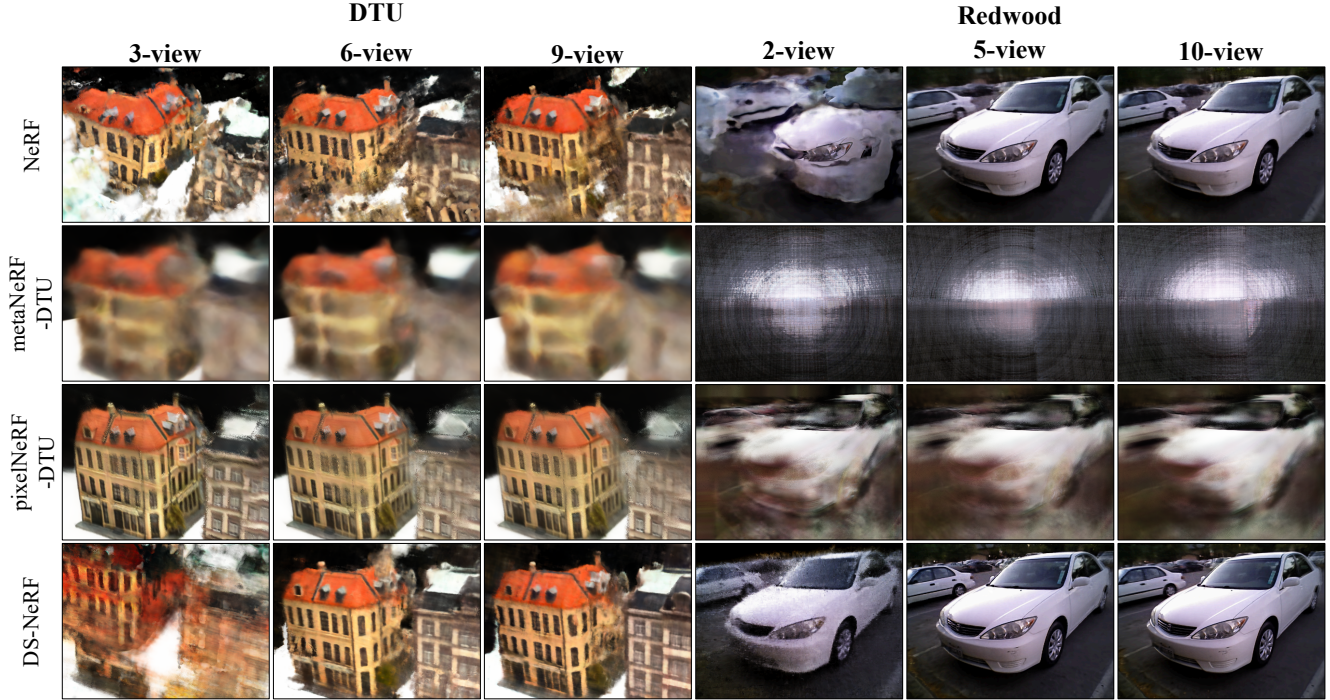Finally we observe mixed results from ablating the confi-

Figure 4: **View Synthesis on DTU and Redwood:** We show view synthesis results for NeRF models given 3, 6, and 9 input views. pixelNeRF, which is pre-trained on DTU, performs the best when given 3-views, although we find DS-NeRF to be visually competitive when more views are available. On Redwood, DS-NeRF is the only baseline to perform well on 2-views.
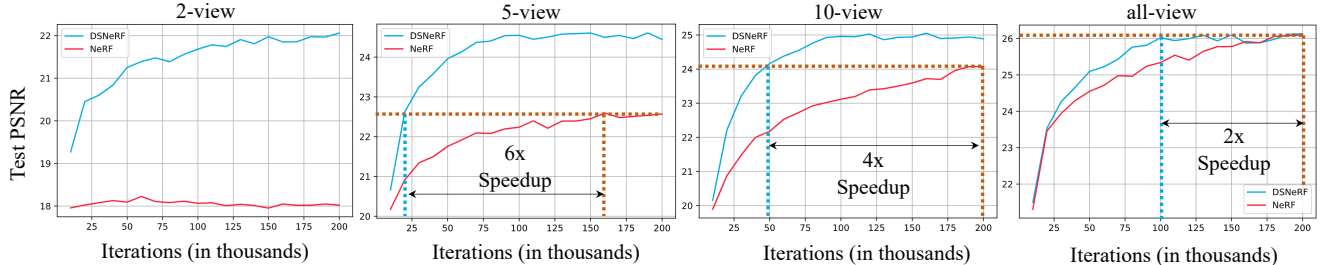


Figure 5: **Faster Training:** We train DS-NeRF and NeRF under identical conditions and observe that DS-NeRF is able to reach NeRF's peak PSNR quality in a fraction of the number of iterations across. On a single RTX 2080Ti, each training iteration of DS-NeRF takes $\sim$ 192.3 ms/iter, while vanilla NeRF needs $\sim$ 191.7 ms/iter. The benefits of depth supervision grows when fewer input views are provided.

dence weighed loss. On NeRF Real, using the weighted loss improves performance, but these benefits disappear when evaluated on DTU. One explanation is that NeRF Real point clouds are more likely to be noisy because they come from a narrow baseline dataset where SFM is less reliable.

## 4.4. Depth error

We can evaluate baseline NeRF's rendered depth by comparing them to reference "ground truth" depth measurements. For NeRF Real dataset, we collect reference depth from test views by running an additional global dense stereo reconstruction. For Redwood-3dscan dataset [3] with released 3d models, we align them with our cameras by running

3dMatch [38], and generate reference depths for each test view. For each baseline, we measure depth error by comparing rendered depth to reference depths. Due to the scale differences between reference and rendered depths, we compute an affine transformation to scale and shift the rendered depth values. To deal with COLMAP's scale ambiguity between different scenes, we report depth error as percentage relative error of the reference depth. Please refer to Appendix A.5 for more details regarding depth error evaluation.

As shown in Table 2, DS-NeRF, trained with supervision obtained only from depth in training views, is able to estimate depth more accurately than all the other NeRF models. While this is not particularly surprising, it does highlight the

8

| Redwood-3dscan | PSNR↑ | | | SSIM↑ | | | LPIPS↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view |
| NeRF | 10.45 | 22.37 | 23.42 | 0.38 | 0.75 | 0.82 | 0.51 | 0.45 | 0.45 |
| metaNeRF-DTU | 14.30 | 14.63 | 15.1 | 0.37 | 0.39 | 0.40 | 0.76 | 0.76 | 0.75 |
| pixelNeRF-DTU | 12.70 | 12.86 | 12.84 | 0.43 | 0.47 | 0.50 | 0.76 | 0.75 | 0.70 |
| DS-NeRF | 18.32 | 22.95 | 23.77 | 0.62 | **0.78** | 0.81 | 0.40 | **0.34** | 0.42 |
| DS-NeRF w/ RGB-D | **22.42** | **23.35** | **23.86** | **0.73** | 0.77 | **0.84** | **0.36** | 0.35 | **0.28** |

Table 5: **View Synthesis on Redwood:** We evaluate view synthesis on 2, 5, and 10 input views on the Redwood dataset [3]. DS-NeRF (with COLMAP[25] inputs) outperforms baselines on various metrics across varying numbers of views. Learning DS-NeRF with the RGB-D reconstruction output [38] further improves performance, highlighting the potential of applying our method to other sources of depth.
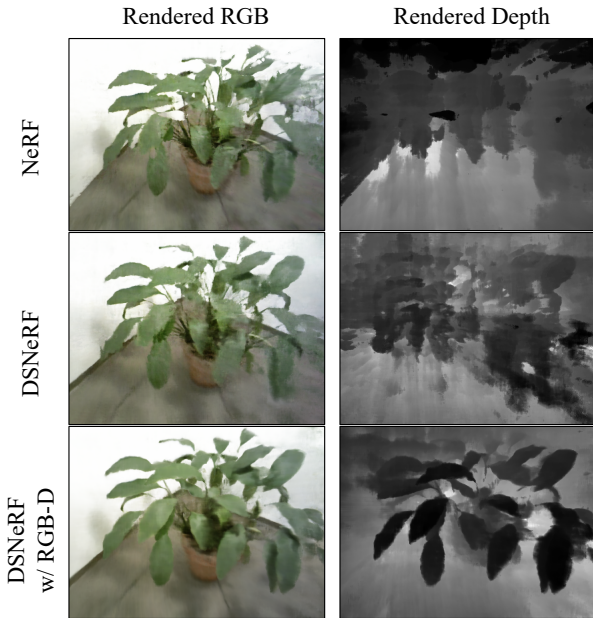
Rendered RGB        Rendered Depth



Figure 6: **Other sources of depth supervision:** Besides COLMAP depth, DS-NeRF is able to use other sources of depth for supervision like an RGB-D input. We derive dense depth maps for each training view using 3DMatch [38] with RGB-D data collected using a depth camera from the Redwood dataset [3] collected by a depth camera. With dense depth supervision, DS-NeRF can render even higher quality images and depth maps than only using structure-from-motion depth keypoints.

weakness of training NeRFs only using RGB supervision. For example, in Figure 3, NeRF tends to ignore geometry and fails to produce any coherent depth map.

**RBG-D inputs.** We consider a variant of depth supervision using RGB-D input from Redwood. Using the procedure described in Appendix A.4, we incorporate significantly more depth supervision. We derive dense depth map for each training view using 3DMatch [38] with RGB-D input. With dense depth supervision, we can now render rays for any pixel in the valid region. As shown in Table 5, Table 2 and

Figure 6, such supervision allows DS-NeRF to render higher quality images from Redwood. While scanned depth may not always be available, there is a current technological trend to integrate such sensors into everyday devices.

### 4.5. Analysis

**Overfitting.** Figure 2 shows that NeRF can overfit to a small number of input views by learning degenerated 3D geometries. Adding depth supervision can disambiguate geometry, and render better novel views.

**Faster Training.** To quantify speed improvements in NeRF training, we compare training DS-NeRF and NeRF under identical settings. Like in Section 4.3, we evaluate view synthesis quality on test views under various number of input views from NeRF Real using PSNR. We can compare training speed performance by plotting PSNR on test views versus training iterations in Figure 5.

DS-NeRF achieves a particular test PSNR threshold using 2-6x less training iterations than NeRF. These benefits are significantly magnified when given fewer views. In the extreme case of only 2-views, NeRF is completely unable to match DS-NeRF's performance and fails. While these results are given in terms of training iteration, we can translate them into wall time improvements. On a single RTX 2080Ti, each training iteration of DS-NeRF takes ∼ 192.3 ms/iter, while vanilla NeRF needs ∼ 191.7 ms/iter. Thus we can see that in the 5-view case, DS-NeRF achieves NeRF's peak test PSNR around 9 hours faster, a massive improvement considering the the negligible cost of depth supervision.

### 5. Discussion

We introduce Depth-supervised NeRF, a model for learning neural radiance fields that takes advantage of depth supervision. Our model uses "free" supervision provided by sparse 3D point clouds computed during standard SFM preprocessing steps. This additional supervision has a significant impact; DS-NeRF trains 2-6X faster and produces better results from fewer training views (improving PSNR from 18.2 to 22.3). While recent research has sought to improve

NeRF by exploiting priors learned from category-specific training data, our approach requires no training and thus generalizes (in principle) to many scenes on which SFM succeeds. Finally, we provide cursory experiments that explore alternate forms of depth supervision such as active depth sensors.

In the future, we are interested in incorporating other types of 3D information into neural radiance field training like outputs from monocular depth networks or video depth and dynamics. Even more interesting would be approaches that can directly initialize high-quality NeRF function parameters from explicit 3D information like point cloud, voxels, and meshes. Such a method would combine the high quality rendering of volumetric representations with the ease and accessibility of explicit 3D representations.

# References

[1] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[2] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[3] Sungjoon Choi, Qian-Yi Zhou, Stephen Miller, and Vladlen Koltun. A large dataset of object scans. *arXiv:1602.02481*, 2016. 2, 5, 6, 8, 9

[4] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, 2017. 3

[5] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Styles Overbeck, Noah Snavely, and Richard Tucker. Deepview: High-quality view synthesis by learned gradient descent. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[6] Ronghang Hu and Deepak Pathak. Worldsheet: Wrapping the world in a 3d sheet for view synthesis from a single image. *arXiv preprint arXiv:2012.09854*, 2020. 3

[7] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 406–413, 2014. 5, 6

[8] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 3

[9] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[10] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3

[11] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. *arXiv preprint arXiv:2105.06466*, 2021. 2

[12] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth estimation. *ACM SIGGRAPH*, 39(4), 2020. 3

[13] Arun Mallya, Ting-Chun Wang, Karan Sapra, and Ming-Yu Liu. World-consistent video-to-video synthesis. In *European Conference on Computer Vision (ECCV)*, 2020. 3

[14] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[15] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[16] Moustafa Meshry, Dan B. Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla. Neural rerendering in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[17] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 2019. 3, 5

[18] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*, 2020. 2, 3, 5, 12

[19] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3

[20] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 2

[21] Thomas Porter and Tom Duff. Compositing digital images. In *Proceedings of the 11th annual conference on Computer graphics and interactive techniques*, pages 253–259, 1984. 3

[22] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[23] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular

depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 3

[24] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 2, 3

[25] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 3, 4, 5, 6, 9

[26] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[27] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3, 4

[28] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2, 3, 4

[29] Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[30] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P. Srinivasan, Jonathan T. Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3, 5, 6, 12

[31] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3, 4

[32] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 5, 6

[33] 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. 5

[34] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. SynSin: End-to-end view synthesis from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3, 4

[35] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. *arXiv preprint arXiv:2011.12950*, 2020. 2

[36] Lin Yen-Chen, Pete Florence, Jonathan T. Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. inerf: Inverting neural radiance fields for pose estimation. *arxiv arXiv:2012.05877*, 2020. 2

[37] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images, 2020. 2, 3, 5, 6, 7, 12

[38] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 8, 9

[39] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5

[40] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *ACM SIGGRAPH*, 2018. 3

# Appendix

We provide additional implementation details and experimental results. Our code and data are available at our GitHub repo and website.

## A. Implementation Details

### A.1. Depth supervision implementation

Depth supervision is implemented by projecting a ray with direction (in local camera coordinates) given by the image coordinates of a detected keypoint and $-1$ in the camera axis. We shoot this ray into a scene and render its depth using the same sampling procedure described in NeRF.

### A.2. COLMAP details

We run the COLMAP with the default configuration on the limited views (the same as NeRF training inputs) (e.g., 2 views). The SfM output is used only during training and is not used during inference. We also normalize the depths for a consistent magnitude of our depth supervision loss.

### A.3. metaNeRF and pixelNeRF baselines

**metaNeRF.** We use the metaNeRF implementation released by Tancik *et al*. [30] which uses Jax. For meta-initialization comparison, we adopt their ShapeNet experimental setup of training a category-specific initialization and further fine-tuning for a fixed number of iterations. In our case, we treat the entire DTU dataset as a single category. We used 64 inner loop optimizations and trained for 40K outer loop steps. We subsequently fine-tune for 1K steps.

For adapting metaNeRF across different domains, we have to deal with different ray bounds and coordinate scales. This is challenging, and we do not directly address this issue as devising approaches to transferring NeRF meta-initialization across different datasets is beyond the scope of this paper. Instead, we used the default scaling and ray bounds provided by previous NeRF works for DTU and NeRF Real [37, 18].

**pixelNeRF-DTU with finetuning.** We start with Yu *et al*. [37]'s pre-trained models on DTU dataset. For a new input scene, we finetune the weights of pixelNeRF on training views for 20K iterations before subsequently evaluating view synthesis on test views.

**pixelNeRF-DTU with finetuning and DS.** This variation is implemented similarly as *pixelNeRF-DTU with finetuning*, but the pixelNeRF finetuning stage incorporates the depth supervision loss.

### A.4. Dataset Splits

**NeRF Real-world.** We split each scene from NeRF Real into training views and test views. Because the number of
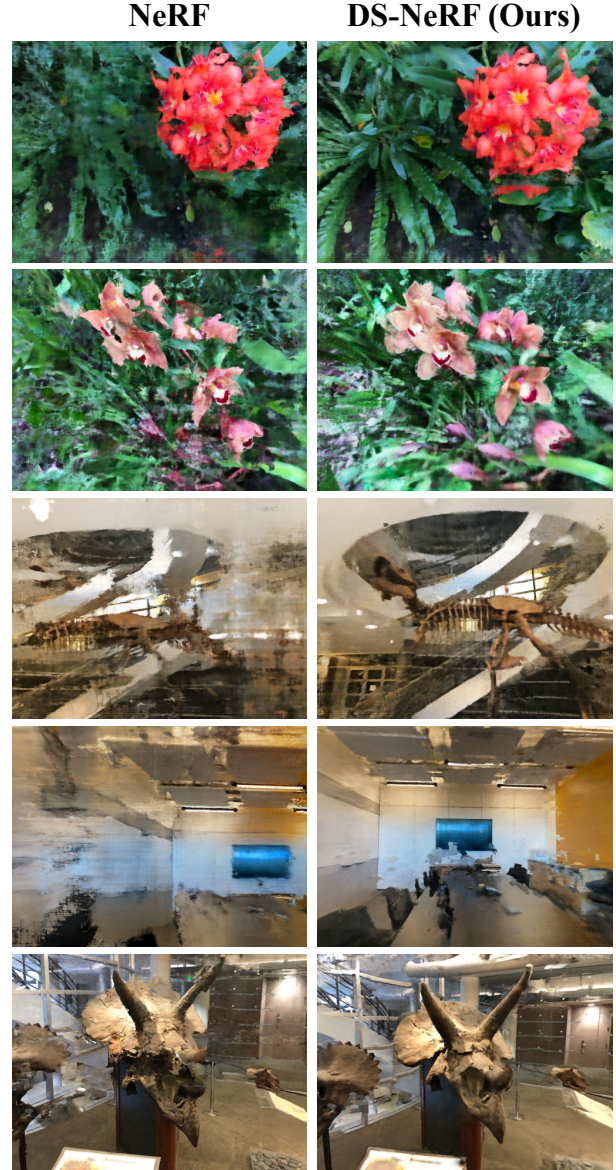
**NeRF**     **DS-NeRF (Ours)**



Figure 7: **2 input view synthesis on NeRF Real:** We render images from a test view for NeRF and DS-NeRF trained on 2 views. DS-NeRF captures better geometry and thus achieves better visual quality when rendering new views. Please see our video for additional insight into the structure of the learnt geometry.

views of each scene varies, we split every eighth image id into the test set (0, 8, 16, 24, . . . ) and construct training views that are evenly distributed over the remaining viewpoint id numbers. This setting gives us sufficient coverage to train different NeRF experiments. We create subsets from these training views of specific sizes to evaluate performance on few-input view synthesis.

**Redwood 3d-scan.** We selected five test scene from the Redwood 3dscan dataset: `table`, `plant`, `chair`, `car`,

| NeRF real-world | PSNR↑ | | | SSIM↑ | | | LPIPS↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view | 2-view | 5-view | 10-view |
| NeRF | 16.31 | 18.37 | 22.29 | 0.48 | 0.49 | 0.53 | 0.55 | 0.51 | 0.53 |
| metaNeRF-NeRF Real | 10.37 | 10.40 | 10.40 | 0.34 | 0.34 | 0.34 | 0.92 | 0.92 | 0.92 |
| pixelNeRF-NeRF Real | 11.56 | 11.59 | 11.61 | 0.36 | 0.37 | 0.37 | 0.87 | 0.87 | 0.86 |
| DS-NeRF | **22.26** | **24.47** | **25.15** | **0.68** | **0.69** | **0.69** | **0.38** | **0.37** | **0.37** |

Table 6: **View Synthesis on NeRF Real (4-fold cross-validation):** Because metaNeRF and pixelNeRF may perform better by training on sequences similar to the validation set, we perform 4-fold cross validation on 8 NeRF Real scenes. We find that training on NeRF Real scenes does not improve metaNerf/pixelNerf performance over training on DTU (c.f. Table 1), possibly because DTU-train has 88 scenes while a 4-fold NeRF Real training dataset has 6 scenes.

and `stool`.

Each scene is constructed using 15 frames of RGB and depth. These 15 frames are further sub-divided into training views and test views. We construct training sets with the following viewpoints [5, 11, 2, 8, 14, 1, 4, 7, 10, 13], truncating when working with a smaller subset (e.g. 2-views uses only 5 and 11). We evaluate the quality of view synthesis on the test views [0, 3, 6, 9, 12].

### A.5. Depth Error Evaluation

To evaluate the depth error of these different baselines, we need to first compute a reference depth of an input scene from test camera poses.

**Depth evaluation on NeRF Real.** For a given set of training views from a scene, we use COLMAP's SfM algorithm to get sparse keypoints and camera poses. In addition, we run dense MVS on all training and test views to get a reference depth map from every view and test poses. To align dense MVS depth with the SfM keypoint depth obtained from the training views, we compute a scale $a$ and shift $b$ scalar that aligns the keypoint depth visible in a training view to its MVS depth. More specifically, we solve the following least squares optimization on those detected keypoints:

$$\min_{a,b} \sum_{p \in \mathcal{P}} (aD_{\text{SfM}}(p) - b - D_{\text{MVS}}(p))^2, \qquad (7)$$

where $\mathcal{P}$ is the set of detected and visible keypoints $p$ for a scene, $D_{\text{SfM}}$ is the sparse depth we use for training, and $D_{\text{MVS}}$ is the dense depth maps we use for evaluation.

Depth error can be computed by transforming the rendered depth $\hat{D}$ from a test camera $c$.

$$\text{Err(D)} = \left\| a\hat{D} - b - D_{\text{MVS}} \right\|^2 \qquad (8)$$

Note that dense MVS depth is only used for depth evaluation. It is not used during training and test by any method.

### B. NeRF Real with pixelNeRF and metaNeRF

While pixelNeRF and metaNeRF can achieve reasonable results when given sufficient pre-training data and a small train-test domain gap, many real-world applications cannot rely on the assumptions of similar test domain and sufficient training samples. We further highlight this by showing what happens to these baselines in such a scenario. To construct a more realistic setting for NeRF-based applications, we perform a 4-fold evaluation by splitting NeRF Real into 6 training scenes and 2 test scenes. The 4 test-splits are: fern-horn, flower-fortress, leaves-orchids, and room-trex.

For every split, we use the remaining 6 training scenes to learn NeRF Real priors for pixelNeRF and metaNeRF and evaluate the view synthesis results after fine-tuning the model on the test scenes. We show these baseline results in Table 6 and find that meta-learning NeRF baselines struggle to properly leverage the priors observed during training.

For metaNeRF, we use the same ShapeNet setup described in Section A.3. We train on NeRF Real training scenes for 40K steps and then finetune for 1K steps on test scenes (*metaNeRF-NeRF Real*). We also tried a slightly longer finetuning duration of 5K steps, but found no significant improvement in results. We observed the following failure mode; within NeRF Real training, certain scenes were dominating the learned prior such that even with many iterations of fine-tuning on test-scenes, NeRF would only render cloudy looking textures from training scenes like `flowers` and `room`.

For pixelNeRF we train the network from scratch on the 6 training scenes. Our training procedure uses 3 input views for each training iteration to finetune the encoder and decoder over 40K iterations. Here we again find that the category-specific models struggle due to the lack of sufficient training data to learn these priors.

### C. Additional Results

We present additional qualitative results of rendered images in Figure 7. We can see that depth supervision significantly improves the quality of test views.