

IBRNet: Learning Multi-View Image-Based Rendering

Qianqian Wang^{1,2} Zhicheng Wang¹ Kyle Genova^{1,3} Pratul Srinivasan¹ Howard Zhou¹
Jonathan T. Barron¹ Ricardo Martin-Brualla¹ Noah Snavely^{1,2} Thomas Funkhouser^{1,3}

¹Google Research ²Cornell Tech, Cornell University ³Princeton University

Abstract

We present a method that synthesizes novel views of complex scenes by interpolating a sparse set of nearby views. The core of our method is a network architecture that includes a multilayer perceptron and a ray transformer that estimates radiance and volume density at continuous 5D locations (3D spatial locations and 2D viewing directions), drawing appearance information on the fly from multiple source views. By drawing on source views at render time, our method hearkens back to classic work on image-based rendering (IBR), and allows us to render high-resolution imagery. Unlike neural scene representation work that optimizes per-scene functions for rendering, we learn a generic view interpolation function that generalizes to novel scenes. We render images using classic volume rendering, which is fully differentiable and allows us to train using only multi-view posed images as supervision. Experiments show that our method outperforms recent novel view synthesis methods that also seek to generalize to novel scenes. Further, if fine-tuned on each scene, our method is competitive with state-of-the-art single-scene neural rendering methods.¹

1. Introduction

Given a set of posed images of a scene, the goal of novel view synthesis is to produce photo-realistic images of the same scene at novel viewpoints. Early work on novel view synthesis focused on image-based rendering (IBR). Starting from the pioneering work on view interpolation of Chen and Williams [6], and proceeding through light field rendering [3, 15, 29], view-dependent texturing [8], and more modern learning-based methods [17], IBR methods generally operate by warping, resampling, and/or blending source views to target viewpoints. Such methods can allow for high-resolution rendering, but generally require either very dense input views or explicit proxy geometry, which is difficult to estimate with high quality leading to artifacts in rendering.

More recently, one of the most promising research directions for novel view synthesis is neural scene representations, which represent scenes as the weights of neural networks. This research area has seen significant progress through the use of Neural Radiance Fields (NeRF) [42]. NeRF shows that multi-layer perceptrons (MLPs) combined with positional encoding can be used to represent the continuous 5D radiance field of a scene, enabling photo-realistic novel view synthesis on complex real-world scenes. NeRF’s use of continuous scene modeling via MLPs, as opposed to explicit discretized volumes [59] or multi-plane images [12, 74] allows for more compact representations and scales to larger viewing volumes.

Although neural scene representations like NeRF can represent scenes faithfully and compactly, they typically require a lengthy optimization process for each new scene before they can synthesize any novel views of that scene, which limits the value of these methods for many real-world applications.

In this work, we leverage ideas from both IBR and NeRF into a new learning-based method that generates a *continuous* scene radiance field on-the-fly from multiple source views for rendering novel views. We learn a general view interpolation function that simultaneously performs density/occlusion/visibility reasoning and color blending while rendering a ray. This enables our system to operate without any scene-specific optimization or precomputed proxy geometry.

At the core of our method is a lightweight MLP network that we call *IBRNet*, which aggregates information from source views along a given ray to compute its final color. For sampled 3D locations along the ray, the network first fetches latent 2D features, derived from nearby source views, that encode spatial context. IBRNet then aggregates these 2D features for each sampled location to produce a *density feature* that captures information about whether that feature seems to be on a surface. A ray transformer module then computes a scalar density value for each sample by considering these density features along the entire ray, enabling visibility reasoning across larger spatial scales. Separately, a

¹<https://ibrnet.github.io/>

color blending module uses the 2D features and view direction vectors from source views to derive a view-dependent color for each sample, computed as a weighted combination of the projected colors of the source views. A final color value is then computed for each ray using volume rendering.

Our approach is fully differentiable and can therefore be trained end-to-end using multi-view images. Our experiments show that when trained on large amounts of data, our method can render high-resolution photo-realistic novel views for unseen scenes that contain complex geometry and materials, and our quantitative evaluation shows that it improves upon state-of-the-art novel view synthesis methods designed to generalize in a single shot to new test scenes. Moreover, for a particular scene, we can fine-tune IBRNet to improve the quality of synthesized novel views to match the performance of state-of-the-art neural scene representation methods like NeRF [42]. In summary, our contributions are:

- a new learning-based multi-view image-based rendering approach that outperforms existing one-shot view synthesis methods on novel scenes,
- a new model architecture called IBRNet that enables the continuous prediction of colors and densities in space from multiple views,
- a per-scene fine-tuning procedure that achieves comparable performance to state-of-the-art novel view synthesis methods designed only for single-scene inference.

2. Related work

Image based rendering. Early work on IBR introduced the idea of synthesizing novel views from a set of reference images by a weighted blending of reference pixels [9, 15, 29]. Blending weights were computed based on ray-space proximity [29] or approximate proxy geometry [3, 9, 19]. In more recent work, researchers have proposed improved methods for computing proxy geometry [5, 18], optical flow correction [4, 10, 11], and soft blending [49, 54]. For example, Hedman et al. [17] use two types of multi-view stereo [23, 56] to produce a view-dependent mesh surface, then use a CNN to compute blending weights. Others synthesize a radiance field directly on a mesh surface [8, 21, 65] or point cloud [1, 40, 50]. While these methods can handle sparser views than other approaches and achieve promising results in some cases, they are fundamentally limited by the performance of 3D reconstruction algorithms [23, 56]. They have difficulty in low-textured or reflective regions, where stereo reconstruction tends to fail, and cannot handle partially translucent surfaces. In contrast, our method learns continuous volume densities in an end-to-end manner that is optimized for synthesis quality, leading to better performance in challenging scenarios.

Volumetric Representations. Another line of work uses discrete volumetric representations to achieve photo-realistic

rendering. Recent methods leverage convolutional neural networks (CNNs) to predict volumetric representations stored in voxel grids [20, 26, 27, 49, 67] or multi-plane images (MPIs) [12, 13, 32, 41, 62, 74]. At test time, novel views can be rendered from these representations via alpha compositing [51]. Trained end-to-end on large datasets, these methods often learn to compensate for the discretization artifacts of low-resolution voxel grids and can generalize reasonably well to different scenes. While they achieve high-quality view synthesis results, they must explicitly process and store large numbers of samples resulting in extensive memory overhead, limiting the resolution of their outputs. In contrast, our method allows for querying color and opacity at continuous 3D locations and 2D viewing directions without storing a full scene representation, and can scale to render high-resolution images. Our method can also handle larger viewing volumes than MPI-based methods.

Neural scene representations. A promising recent research direction is the use of neural networks for representing the shape and appearance of scenes. Earlier work [2, 14, 24, 39, 44, 47, 70] showed that MLPs can be used as implicit shape representations, where the weights of the MLPs map continuous spatial coordinates to signed distance or occupancy values. With advances in differentiable rendering methods [7, 25, 30, 34, 35, 43], many methods [33, 36, 43, 55, 57, 59, 60, 64, 71] showed the ability to learn both scene geometry and appearance from multi-view observations for simple geometry and diffuse materials. The more recent NeRF [42] achieves very impressive results for novel view synthesis by optimizing a 5D neural radiance field for a scene, suggesting the advantages of continuous representations over discrete ones like voxel grids or meshes. While NeRF opens up many new research opportunities [31, 38, 45, 48, 58, 61], it must be optimized for each new scene, taking hours or days to converge. Concurrent work [66, 72] tries to address this issue and generalize NeRF (with a focus on very sparse input views). However, their use of absolute locations as direct network inputs restricts their ability to generalize to arbitrary new scenes. In contrast, our method can generalize to new, real-world scenes with high quality.

3. Method

Given nearby source views, our method uses volume rendering to synthesize a target view at a novel camera pose. The core problem we try to solve is to obtain colors and densities in a continuous space by aggregating information present in the source views. Our system pipeline (Fig. 1) can be divided into three parts: 1) identifying a set of nearby source views as input and extracting dense features from each source view, 2) predicting volume densities σ and colors c at continuous 5D locations (3D spatial locations and

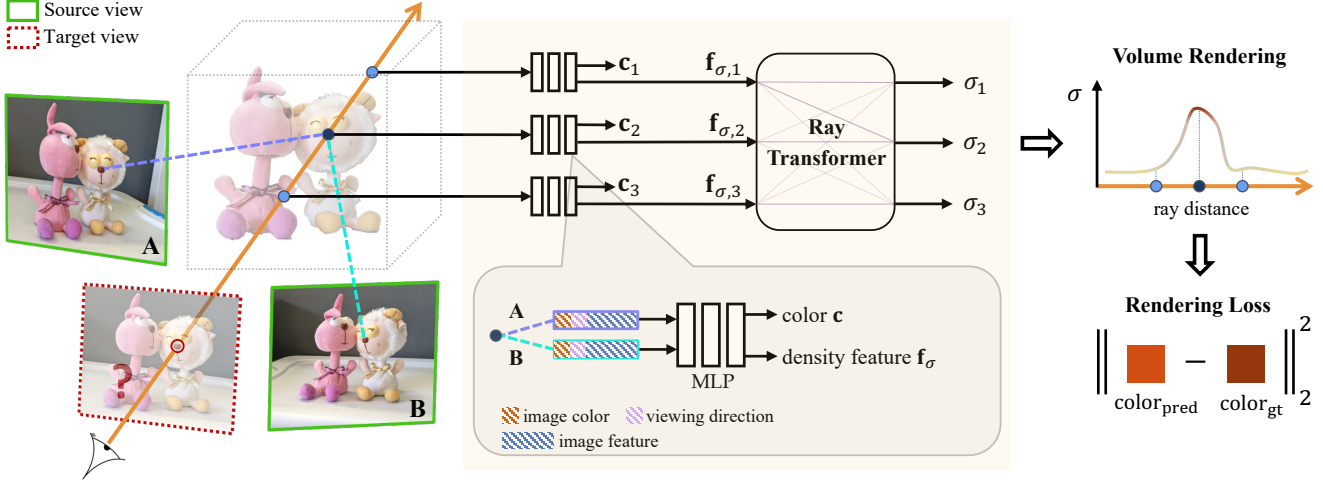


Figure 1: **System Overview.** 1) To render a novel target view (shown here as the image labeled with a ‘?’), we first identify a set of neighboring source views (e.g., the views labeled **A** and **B**) and extract their image features. 2) Then, for each ray in the target view, we compute colors and densities for a set of samples along the ray using our proposed IBRNet (yellow shaded region). Specifically, for each sample, we aggregate its corresponding information (image colors, features, and viewing directions) from the neighboring source views to produce its color \mathbf{c} and density features \mathbf{f}_σ (note that these features are not yet scalar density values). We then apply our proposed ray transformer to these density features across all samples on the ray to predict densities. 3) Finally, we use volume rendering to accumulate colors and densities along the ray to render its color. Our method can be trained end-to-end with an L_2 loss on reconstructed image colors.

2D view directions), and 3) compositing those colors and densities along each camera ray through volume rendering to produce a synthesized image.

3.1. View selection and feature extraction

Unlike neural scene representations that attempt to encode an entire scene into a single network, we synthesize the novel target view by interpolating nearby source views. While our network (described in the next section) can handle an arbitrary number of neighboring views, given limited GPU memory we select a small number of source views as the “working set” for rendering a novel view. To obtain an effective working set, we identify spatially nearby candidate views, then select the subset of N views whose viewing directions are most similar to the target view.

Let $\mathbf{I}_i \in [0, 1]^{H_i \times W_i \times 3}$ and $\mathbf{P}_i \in \mathbb{R}^{3 \times 4}$ respectively denote the color image and camera projection matrix for the i -th source view. We use a shared U-Net based convolutional neural network to extract dense features $\mathbf{F}_i \in \mathbb{R}^{H_i \times W_i \times d}$ from each image \mathbf{I}_i . The set of tuples $\{(\mathbf{I}_i, \mathbf{P}_i, \mathbf{F}_i)\}_{i=1}^N$ forms the input for rendering a target view.

3.2. RGB- σ prediction using IBRNet

Our method synthesizes images based on classic volume rendering, accumulating colors and densities in the 3D scene to render 2D images. We propose IBRNet (illustrated in Fig. 2) to predict colors and densities at continuous 5D locations by aggregating information from multiple source views and incorporating long-range context along the ray.

Our proposed IBRNet is permutation-invariant and accepts a variable number of source views.

We first describe the input to IBRNet for a single query point location $\mathbf{x} \in \mathbb{R}^3$ on a ray $\mathbf{r} \in \mathbb{R}^3$ with unit-length viewing direction $\mathbf{d} \in \mathbb{R}^3$. We project \mathbf{x} into all source views using their camera parameters, and extract colors and features at the projected pixel locations through bilinear interpolation. Let $\{C_i\}_{i=1}^N \in [0, 1]^3$ and $\{\mathbf{f}_i\}_{i=1}^N \in \mathbb{R}^d$ respectively denote these extracted colors and image features for point \mathbf{x} . We also take into account the viewing directions for \mathbf{x} in all source views, denoted as $\{\mathbf{d}_i\}_{i=1}^N$.

3.2.1 Volume density prediction

Our density prediction at point (\mathbf{x}, \mathbf{d}) involves two steps. First, we aggregate the multi-view features at (\mathbf{x}, \mathbf{d}) to obtain a density feature. Then, our proposed ray transformer takes in density features for all samples on a ray and incorporates long-range contextual information to predict the density for each sample including (\mathbf{x}, \mathbf{d}) .

Multi-view feature aggregation. We observe that 3D points on surfaces are more likely to have consistent local appearances in multiple views than 3D points in free space. Therefore, an effective way to infer density would be to check the consistency among the features $\{\mathbf{f}_i\}_{i=1}^N$ for a given point, and a possible implementation would be a PointNet-like [52] architecture that takes in multi-view features and uses variance as the global pooling operator. Specifically, we first compute a per-element mean $\boldsymbol{\mu} \in \mathbb{R}^d$ and

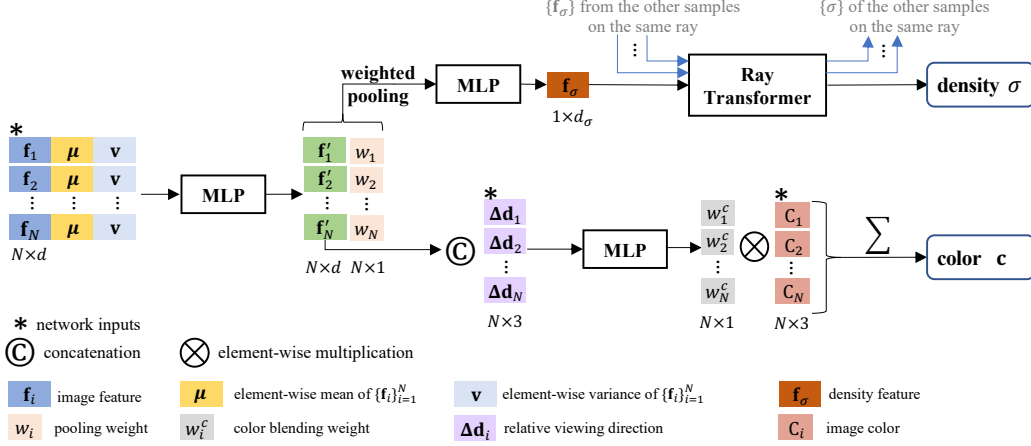


Figure 2: **IBRNet for volume density and color prediction at a continuous 5D location (x, d) .** We first input the 2D image features $\{f_i\}_{i=1}^N$ extracted from all source views to a PointNet-like MLP to aggregate local and global information, resulting in multi-view aware features $\{f'_i\}_{i=1}^N$ and pooling weights $\{w_i\}_{i=1}^N$. To predict density, we pool $\{f'_i\}_{i=1}^N$ using weights $\{w_i\}_{i=1}^N$ which enables multi-view visibility reasoning to obtain a density feature f_σ . Instead of directly predicting density σ from f_σ for individual 5D samples, we use a *ray transformer* module to aggregate information of all samples along the ray. The ray transformer module takes f_σ for all samples on a ray and predicts all their densities (only the density output for (x, d) is highlighted in the figure for simplicity). The ray transformer module enables geometric reasoning across a longer range and improves density predictions. For color prediction, we concatenate $\{f'_i\}_{i=1}^N$ with the viewing directions of the query ray relative to each viewing direction of the source view, i.e., $\{\Delta d_i\}_{i=1}^N$, and predict a set of blending weights. The output color c is a weighted average of the image colors from the source views.

variance $v \in \mathbb{R}^d$ from features $\{f_i\}_{i=1}^N$ to capture global information, and concatenate each f_i with μ and v . Each concatenated feature is fed into a small shared MLP to integrate both local and global information, resulting in a multi-view aware feature f'_i and a weight vector $w_i \in [0, 1]$. We pool these new features $\{f'_i\}_{i=1}^N$ by computing their weighted average and variance using weights $\{w_i\}_{i=1}^N$, which are mapped to a density feature $f_\sigma \in \mathbb{R}^{d_\sigma}$ using an MLP. Compared to a direct average or max-pooling in a PointNet [52], we find that our weighted pooling improves the network’s ability to handle occlusions.

Ray transformer. After obtaining a density feature f_σ , one could directly turn it into a single density σ with another MLP. However, we find such an approach fails to predict accurate densities for new scenes with complex geometry. We ascribe this to the fact that looking at features for a single point *in isolation* is inadequate for accurately predicting its density, and more contextual information is needed—similar to how plane-sweep stereo methods consider matching scores along a whole ray before determining the depth of a particular pixel.

We therefore introduce a new ray transformer module to enable samples on a ray to attend to each other before predicting their densities. The ray transformer consists of the two core components of a classic Transformer [69]: positional encoding and self-attention. Given M samples along the ray, our ray transformer treats samples from near to far as a sequence, and applies positional encoding and

multi-head self-attention to the sequence of density features $(f_\sigma(x_1), \dots, f_\sigma(x_M))$. The final density value σ for each sample is then predicted from its attended feature. The ray transformer module allows for a variable number of inputs and only introduces only a small increase in parameters and computational overhead, while dramatically improving the quality of the predicted densities and the final synthesized image, as shown in our ablation studies.

Improving temporal visual consistency. Our method considers only nearby source views as the working set when synthesizing a target view. Therefore, when generating videos along smooth camera paths, our method is potentially subject to temporarily inconsistent density predictions and flickering artifacts due to abrupt changes in the working set as the camera moves.

To alleviate this issue, we adopt the pooling technique of Sun *et al.* [63]. We replace the mean μ and variance v of $\{f_i\}_{i=1}^N$ with a weighted mean μ_w and variance v_w , weighted so as to reduce the influence of the furthest images in the working set. The weight function is defined as:

$$\tilde{w}_i^f(d, d_i) = \max\left(0, e^{s(d \cdot d_i - 1)} - \min_{j=1 \dots N} e^{s(d \cdot d_j - 1)}\right),$$

$$w_i^f(d, d_i) = \frac{\tilde{w}_i^f(d, d_i)}{\sum_j \tilde{w}_j^f(d, d_j)}, \quad (1)$$

where s is a learnable parameter. This technique improves our method in two ways: 1) it smooths the change in global

features between adjacent frames, and 2) it produces more reasonable global features by up-weighting closer views in the working set and down-weighting more distant views. We observe this technique empirically improves synthesis stability and quality.

3.2.2 Color prediction

We obtain color at a 5D point by predicting blending weights for the image colors $\{C_i\}_{i=1}^N$ in the source views that correspond to that 5D point. Unlike NeRF [42], which uses absolute viewing directions, we consider viewing direction relative to that of the source views, i.e., the difference between \mathbf{d} and \mathbf{d}_i . A smaller difference between \mathbf{d} and \mathbf{d}_i typically implies a larger chance of the color at target view resembling the corresponding color at view i , and vice versa.

To predict the blending weight for each source color C_i , we concatenate the feature \mathbf{f}_i^c with $\Delta\mathbf{d}_i = \mathbf{d} - \mathbf{d}_i$ and input each concatenated feature into a small network to yield a blending weight w_i^c . The final color for this 5D location is blended via a soft-argmax operator $\mathbf{c} = \sum_{i=1}^N [C_i \exp(w_i^c) / \sum_{j=1}^N \exp(w_j^c)]$. An alternative to predicting blending weights would be to directly regress \mathbf{c} . However, we found that direct regression of \mathbf{c} led to decreased performance.

3.3. Rendering and training

The approach introduced in the previous section allows us to compute the color and density value at continuous 5D location. To render the color of a ray \mathbf{r} passing through the scene, we first query the colors and densities of M samples on the ray, and then accumulate colors along the ray modulated by densities:

$$\tilde{C}(\mathbf{r}) = \sum_{k=1}^M T_k (1 - \exp(-\sigma_k)) \mathbf{c}_k, \quad (2)$$

$$\text{where } T_k = \exp\left(-\sum_{j=1}^{k-1} \sigma_j\right). \quad (3)$$

Here samples from 1 to M are sorted to have ascending depth values. \mathbf{c}_k and σ_k denote the color and density of the k -th sample on the ray, respectively.

Hierarchical volume rendering. One benefit of having continuous RGB- σ predictions is that it allows more adaptive and efficient sampling in space. Following NeRF [42], we perform hierarchical volume sampling and simultaneously optimize two networks, a coarse IBRNet and a fine IBRNet, with identical network architecture. At the coarse scale, we sample a set of M_c locations at equidistant disparity (inverse depth) which results in equal intervals between adjacent point projections in pixel space. Given the prediction of the coarse network, we then conduct a more informed sampling

of points along each ray, where samples are more likely to be located at relevant regions for rendering. We sample an additional M_f locations and use all $M_c + M_f$ locations to render the fine results as in NeRF [42]. In our network, we set both M_c and M_f to 64.

Training objective. We render the color of each ray using both the coarse and fine set of samples, and minimize the mean squared error between the rendered colors and ground-truth pixel colors for training:

$$\mathcal{L} = \sum_{\mathbf{r} \in R} \left[\left\| \tilde{C}_c(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 + \left\| \tilde{C}_f(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 \right] \quad (4)$$

where R is the set of rays in each training batch. This allows us to train the feature extraction network as well as the coarse and fine IBRNet simultaneously.

Fine-tuning. Trained on large amounts of data, our method is able to generalize well to novel scenes. Our method can also be fine-tuned per scene using the same objective at training time (Eq. 4) but only on a specific scene, thereby improving synthesis performance on that scene.

3.4. Implementation details

Source and target view sampling. Given multiple views of a scene, we construct a training pair of source and target view by first randomly selecting a target view, and then sampling N nearby views as source views. To select source views, we first identify a pool of $n \cdot N$ nearby views (n is sampled uniformly at random from $[1, 3]$), and then randomly sample N views from the pool. This sampling strategy simulates various view densities during training and therefore helps the network generalize across view densities. During training, we sample N uniformly at random from $[8, 12]$.

Network details. We implement the image feature extraction network using a U-Net like architecture with ResNet34 [16] truncated after layer3 as the encoder, and two additional up-sampling layers with convolutions and skip-connections as the decoder. We decode two sets of feature maps in the final decoding layer, to be used as input to the coarse and fine IBRNet respectively. Both coarse and fine feature maps have $d = 32$ dimensions, and are $1/4 \times$ the original image size. For IBRNet, the dimension of the density feature \mathbf{f}_σ is 16 and we use 4 heads for the self-attention module in ray transformer. Please refer to the supplementary material for more network details.

Training details. We train both the feature extraction network and the IBRNet end-to-end on datasets of multi-view posed images using Adam [28]. The base learning rates for feature extraction network and IBRNet are 10^{-3} and 5×10^{-4} , respectively, which decay exponentially along with the optimization. During fine-tuning, we optimize both our 2D feature extractor and IBRNet itself using smaller

Method	Settings	Diffuse Synthetic 360° [59]			Realistic Synthetic 360° [42]			Real Forward-Facing [41]		
		PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
LLFF [41]	No per-scene	34.38	0.985	0.048	24.88	0.911	0.114	24.13	0.798	0.212
Ours	optimization	37.17	0.990	0.017	25.49	0.916	0.100	25.13	0.817	0.205
SRN [60]	Per-scene optimization	33.20	0.963	0.073	22.26	0.846	0.170	22.84	0.668	0.378
NV [36]		29.62	0.929	0.099	26.05	0.893	0.160	-	-	-
NeRF [42]		40.15	0.991	0.023	31.01	0.947	0.081	26.50	0.811	0.250
Ours _{ft}		42.93	0.997	0.009	28.14	0.942	0.072	26.73	0.851	0.175

Table 1: **Quantitative comparison on datasets of synthetic and real images.** Our evaluation metrics are PSNR/SSIM (higher is better) and LPIPS [73] (lower is better). Both Ours and LLFF [41] are trained on large amounts of training data and then evaluated on all test scenes without any per-scene tuning. Ours consistently outperforms LLFF [41] on all datasets. We also compare our method with neural rendering methods (SRN [60], NV [36], and NeRF [42]) that train a separate network for each scene. To compete fairly with these methods, we also fine-tune our pretrained model on each scene (Ours_{ft}). After fine-tuning, Ours_{ft} is competitive with the state-of-the-art method NeRF [42].

base learning rates (5×10^{-4} and 2×10^{-4}). For pretraining, we train on eight V100 GPUs with a batch size of 3,200 to 9,600 rays depending on image resolution and the number of source views, which takes about a day to finish.

4. Experiments

We evaluate our method in two ways: a) directly evaluating our pretrained model on all test scenes without any fine-tuning (Ours). Note that we train only *one* model and evaluate it on all test scenes. And b) fine-tuning our pretrained model on each test scene before evaluation (Ours_{ft}), similar to NeRF.

4.1. Experimental Settings

Training datasets. Our training data consists of both synthetic data and real data. For synthetic data, we generate object-centric renderings of the 1,023 models in Google Scanned Objects [53]. For real data, we use RealEstate10K [74], the Spaces dataset [12], and 102 real scenes from handheld cellphone captures (35 from LLFF [41] and 67 from ourselves). RealEstate10K is a large video dataset of indoor scenes with camera poses. The Spaces dataset contains 100 scenes captured by a 16-camera rig. Each of the cellphone captured scenes has 20-60 images captured by forward-facing cameras distributed roughly on 2D grids [41]. We use COLMAP [56] to estimate camera parameters and scene bounds for our captures. Our training data includes various camera setups and scene types, which allows our method to generalize well to unseen scenarios.

Evaluation datasets. We use both synthetic rendering of objects and real images of complex scenes for evaluation, as in NeRF [42]. The synthetic data consists of four Lambertian objects with simple geometry from the DeepVoxels [59] dataset, and eight objects with complicated geometry and realistic materials from the NeRF synthetic dataset. Each object in DeepVoxels has 479 training views and 1,000 test

views at 512×512 resolution sampled on the upper hemisphere. Objects in the NeRF synthetic dataset have 100 training views and 200 test views at 800×800 resolution sampled either on the upper hemisphere or the full sphere. The real dataset contains 8 complex real-world scenes captured with roughly forward-facing images from NeRF. Each scene is captured with 20 to 62 images at 1008×756 resolution. 1/8th of the images are held out for testing.

Baselines. We compare our method against the current best view synthesis methods from two classes: *MPI-based methods* and *neural rendering methods*. For MPI-based methods, we compare Ours with a state-of-the-art method LLFF [41] which uses a 3D convolutional neural network to predict a discretized grid for each input image and blends nearby MPIs into the novel viewpoint. Both Ours and LLFF are designed to generalize to novel scenes. We compare Ours_{ft} with neural rendering methods: Neural Volumes (NV) [36], Scene Representation Networks (SRN) [60], and NeRF [42]. These methods train separate networks for each scene or each class of scenes, and either do not generalize at all, or only generalize within object categories. When evaluating against these neural scene representation methods we fine-tune our model per-scene, as do the baselines.

4.2. Results

To render each test view we sample 10 source views from the training set for all evaluation datasets. We evaluate our method and our baselines using PSNR, SSIM, and LPIPS [73]. Results can be seen in Tab. 1 and in Fig. 3.

Tab. 1 shows that our pretrained model generalizes well to novel scenes and consistently outperforms LLFF [41] on all test scenes. After fine-tuning, our model’s performance is competitive with state-of-the-art neural rendering methods like NeRF [42]: We outperform NeRF on both *Diffuse Synthetic 360°* and *Real Forward-Facing* [41] and only have lower PSNR and SSIM on *Realistic Synthetic 360°*. One reason that we underperform NeRF on *Realistic Synthetic 360°*

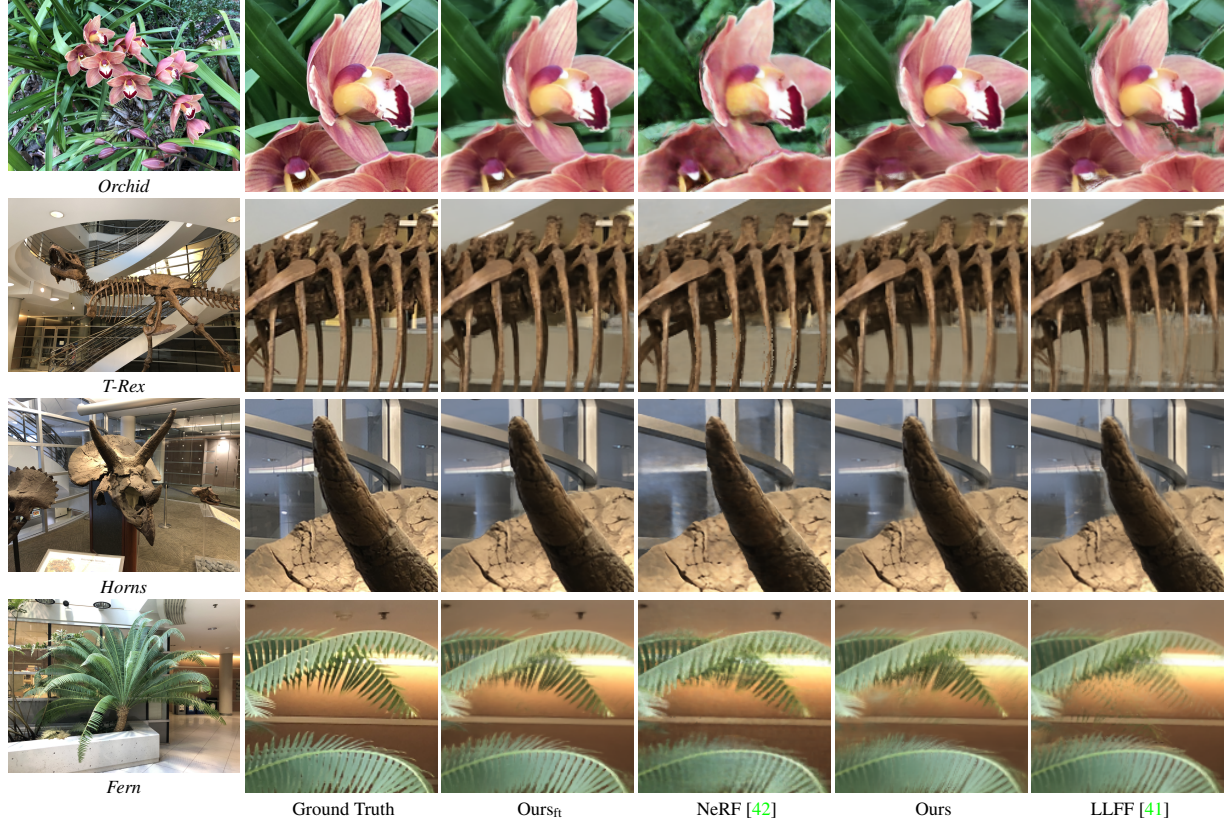


Figure 3: **Qualitative comparison on *Real-Forward-Facing* data.** Our method can more accurately recover fine details in both geometry and appearance, and produce images that are perceptually more similar to ground-truth than other methods. LLFF [41] has difficulty recovering clean and accurate boundary (ghosting artifacts in *Orchid* and repeated edges in *Horns*), and fails to capture thin structures (ribs in *T-Rex*) or partially occluded origins (leaves in *Fern*). NeRF [42]-rendered images exhibits unrealistic noise in *Orchid*. Meanwhile, the texture on the petals is missing and the region near the boundary of the petals are not well recovered. Our method is also slightly better than NeRF at fine structures in *T-Rex* and *Fern*, and can reconstruct more details of the reflection on the glass in *Horns*.

is that this data has very sparse training views for scenes with very complex geometry. Therefore, the information contained in the limited set of local source views may be insufficient for synthesizing a novel view, in which case methods that optimize a global radiance field using all views like NeRF may be better suited.

On the *Real Forward-Facing* [41] data Ours_{fit} achieves substantially better SSIM and LPIPS than NeRF [42], indicating that our synthesized images look more photo-realistic. A key component in NeRF is the use of a positional encoding [64] which helps generate high-frequency details. But positional encoding may also cause unwanted high-frequency artifacts in images (see Fig. 3 *Orchid*), reducing perceptual quality. In contrast, our model does not use positional encoding to regress colors but instead blends colors from source views, biasing our synthesized images to look more like natural images. Our method produces reasonable proxy geometry and can be used to produce temporarily consistent videos. Please see the supplemental material for more detail.

4.3. Ablations and analysis

Ablation studies. We conduct ablation studies (Tab. 2) to investigate the individual contribution of key aspects of our method, using our pretrained models on *Real Forward-Facing* [41] data. For “No ray transformer”, we remove the ray transformer so that the density prediction for samples on each ray is independent and local. Removing this module causes the network to fail to predict accurate densities, leading to “black hole” artifacts and blurriness (see supplementary material for qualitative comparison). For “No view directions”, we remove the view direction input in our network. This includes removing the weighted pooling in Eq. 1 and removing the viewing direction as a network input for color prediction. This reduces our model’s ability to reproduce view-dependent effects such as specularities, but does not reduce performance as significantly as seen in NeRF because our model blends nearby input images (which themselves contain view-dependent effects). For “Direct color regression”, we directly predict the RGB values at each 5D

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
No ray transformer	21.31	0.675	0.355
No view directions	24.20	0.796	0.243
Direct color regression	24.73	0.810	0.220
Full model Ours	25.13	0.817	0.205

Table 2: **Ablation study on Real Forward-Facing [41] data.** We report the metrics for the pretrained model of each ablation without per-scene fine-tuning.

location instead of the blending weights. This is modestly worse than blending image colors from source views.

Sensitivity to source view density. We investigate how our method degrades as the source views become sparser on *Diffuse Synthetic 360°* [59]. Each scene in the original data provides 479 training views randomly distributed on the upper hemisphere. We uniformly sub-sample the training views by factors of 2, 4, 6, 8, 10 to create varying view densities. We report the PSNR for both the pretrained model Ours and the fine-tuned model Ours_{fit} in Fig. 4. Our method degrades reasonably as the input views become sparser.

Computational overhead at inference. Here we investigate the computation required by IBRNet in comparison to other methods. Our inference pipeline for rendering a target view can be divided into two stages: We first extract image features from all source views of interest, which must only be done once as features can then be used repeatedly for rendering. Given these pre-computed features we use our proposed IBRNet to produce the colors and densities for samples along each camera ray. Because image features can be precomputed, the major computational cost at inference is incurred by IBRNet. Since our method interpolates novel views from a set of source views, the #FLOPs of IBRNet depends on the size of the local working set, i.e., the number of source views (#Src.Views). We report the size of IBRNet and #FLOPs required to render a single pixel under varying numbers of source views in Tab. 3.

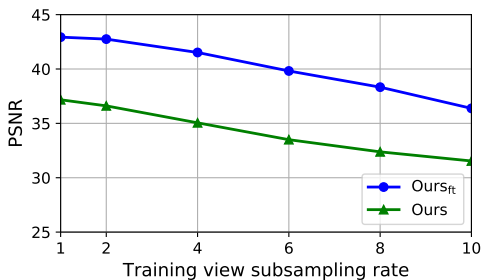


Figure 4: **Sensitivity to source view density.** We subsample the views provided in *Diffuse Synthetic 360°* [59] by {2, 4, 6, 8, 10} to create varying source view densities on the upper hemisphere.

Method	#Params	#Src.Views	#FLOPs	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
SRN	0.55M	-	5M	22.84	0.668	0.378
NeRF	1.19M	-	304M	26.50	0.811	0.250
Ours _{fit}	0.04M	5	29M	25.80	0.828	0.190
		8	45M	26.56	0.847	0.176
		10	55M	26.73	0.851	0.175

Table 3: **Network size and computational cost at inference.** The network size of our MLP is much smaller than SRN [60] and NeRF [42]. All #FLOPs reported are for rendering a single pixel. Both NeRF [42] and Ours_{fit} use hierarchical volume sampling. M_c , M_f are set to 64, 128 respectively for NeRF, and 64, 64 for Ours_{fit}. SRN [60] uses ray marching with only 10 steps thus having much fewer #FLOPs. For our method, the number of #FLOPs scales roughly linearly with the number of source views (#Src.Views).

We see that IBRNet requires less than 20% of the #FLOPs than NeRF requires while achieving comparable performance. We consider two explanations for IBRNet’s increased efficiency. First, NeRF uses a single MLP to encode the whole scene. To query any point in the scene, one must touch every parameter of the MLP. As the scene becomes more complex, NeRF needs more network capacity, which leads to larger #FLOPs for each query. In contrast, our method interpolates nearby views locally, and therefore the #FLOPs does not grow with the scale of the scene but only with the size of the local working set. This property also makes our approach more dynamic, enabling applications like streaming large scenes and adjusting computational costs on-the-fly to suit a device. Second, we encode the scene as posed images and features, a more structured and interpretable representation than the MLP used in NeRF. A simple projection operation allows us to fetch the most relevant information from source views for estimating radiance at a point in space. IBRNet must only learn how to combine that information, not how to synthesize it from scratch.

5. Conclusion

We proposed a learning-based multi-view image-based rendering framework that synthesizes novel views of a scene by blending pixels from nearby images with weights and volume densities inferred by a network comprising an MLP and ray transformer. This approach combines the advantages of IBR and NeRF to produce state-of-the-art rendering quality on complex scenes without requiring precomputed geometry (unlike many IBR methods), storing expensive discretized volumes (unlike neural voxel representations), or expensive training for each new scene (unlike NeRF).

Acknowledgements We thank Jianing Wei, Liangkai Zhang, Adel Ahmadyan, and Bhav Ashok for helpful discussions.

References

- [1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. *arXiv preprint arXiv:1906.08240*, 2019. 2
- [2] Matan Atzmon, Niv Haim, Lior Yariv, Ofer Israelov, Haggai Maron, and Yaron Lipman. Controlling neural level sets. *NeurIPS*, 2019. 2
- [3] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. *SIGGRAPH*, 2001. 1, 2
- [4] Dan Casas, Christian Richardt, John Collomosse, Christian Theobalt, and Adrian Hilton. 4d model flow: Precomputed appearance alignment for real-time 4d video interpolation. *Computer Graphics Forum*, 2015. 2
- [5] Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM TOG*, 2013. 2
- [6] Shenchang Eric Chen and Lance Williams. View interpolation for image synthesis. *SIGGRAPH*, 1993. 1
- [7] Wenzheng Chen, Jun Gao, Huan Ling, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. *NeurIPS*, 2019. 2
- [8] Paul Debevec, Yizhou Yu, and George Borshukov. Efficient view-dependent image-based rendering with projective texture-mapping. *Eurographics Workshop on Rendering Techniques*, 1998. 1, 2
- [9] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. *Computer graphics and interactive techniques*, 1996. 2
- [10] Ruofei Du, Ming Chuang, Wayne Chang, Hugues Hoppe, and Amitabh Varshney. Montage4d: Interactive seamless fusion of multiview video textures. *ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games*, 2018. 2
- [11] Martin Eisemann, Bert De Decker, Marcus Magnor, Philippe Bekaert, Edilson De Aguiar, Naveed Ahmed, Christian Theobalt, and Anita Sellent. Floating textures. *Computer graphics forum*, 2008. 2
- [12] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. *CVPR*, 2019. 1, 2, 6
- [13] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world’s imagery. *CVPR*, 2016. 2
- [14] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. *CVPR*, 2020. 2
- [15] Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. The Lumigraph. *SIGGRAPH*, 1996. 1, 2
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*, 2016. 5, 12
- [17] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *SIGGRAPH Asia*, 2018. 1, 2
- [18] Peter Hedman, Tobias Ritschel, George Drettakis, and Gabriel Brostow. Scalable inside-out image-based rendering. *ACM TOG*, 2016. 2
- [19] Benno Heigl, Reinhard Koch, Marc Pollefeys, Joachim Denzler, and Luc Van Gool. Plenoptic modeling and rendering from image sequences taken by a hand-held camera. 1999. 2
- [20] Philipp Henzler, Niloy J. Mitra, and Tobias Ritschel. Learning a neural 3d texture space from 2d exemplars. *CVPR*, 2020. 2
- [21] Jingwei Huang, Justus Thies, Angela Dai, Abhijit Kundu, Chiyu Jiang, Leonidas J Guibas, Matthias Nießner, and Thomas Funkhouser. Adversarial texture optimization from rgb-d scans. *CVPR*, 2020. 2
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015. 12
- [23] Michal Jancosek and Tomáš Pajdla. Multi-view reconstruction preserving weakly-supported surfaces. *CVPR*, 2011. 2
- [24] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3d scenes. *CVPR*, 2020. 2
- [25] Yue Jiang, Dantong Ji, Zhizhong Han, and Matthias Zwicker. Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. *CVPR*, 2020. 2
- [26] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM TOG*, 2016. 2
- [27] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. *NeurIPS*, 2017. 2
- [28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 5
- [29] Marc Levoy and Pat Hanrahan. Light field rendering. *SIGGRAPH*, 1996. 1, 2
- [30] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 37(6):222:1–222:11, 2018. 2
- [31] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. *arXiv preprint arXiv:2011.13084*, 2020. 2
- [32] Zhengqi Li, Wenqi Xian, Abe Davis, and Noah Snavely. Crowdsampling the plenoptic function. *ECCV*, 2020. 2
- [33] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020. 2
- [34] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. *ICCV*, 2019. 2
- [35] Shaohui Liu, Yinda Zhang, Songyou Peng, Boxin Shi, Marc Pollefeys, and Zhaopeng Cui. Dist: Rendering deep implicit signed distance function with differentiable sphere tracing. *CVPR*, 2020. 2
- [36] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM TOG*, 2019. 2, 6

- [37] Keyang Luo, Tao Guan, Lili Ju, Yuesong Wang, Zhuo Chen, and Yawei Luo. Attention-aware multi-view stereo. *CVPR*, 2020. 12
- [38] 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. *CVPR*, 2021. 2
- [39] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. *CVPR*, 2019. 2
- [40] Moustafa Mahmoud Meshry, Dan B Goldman, Sameh Khamis, Hugues Hoppe, Rohit Kumar Pandey, Noah Snavely, and Ricardo Martin Brualla. Neural rerendering in the wild. *CVPR*, 2019. 2
- [41] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima K. Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM TOG*, 2019. 2, 6, 7, 8, 12
- [42] 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. *ECCV*, 2020. 1, 2, 5, 6, 7, 8, 12, 15
- [43] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. *CVPR*, 2020. 2
- [44] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *CVPR*, 2019. 2
- [45] 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
- [46] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. *NIPS-W*, 2017. 12
- [47] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. *ECCV*, 2020. 2
- [48] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. *CVPR*, 2021. 2
- [49] Eric Penner and Li Zhang. Soft 3D reconstruction for view synthesis. *SIGGRAPH Asia*, 2017. 2
- [50] Francesco Pittaluga, Sanjeev J Koppal, Sing Bing Kang, and Sudipta N Sinha. Revealing scenes by inverting structure from motion reconstructions. *CVPR*, 2019. 2
- [51] Thomas Porter and Tom Duff. Compositing digital images. *Computer graphics and interactive techniques*, 1984. 2
- [52] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CVPR*, 2017. 3, 4
- [53] Google Research. Google scanned objects. <https://app.ignitionrobotics.org/GoogleResearch/fuel/collections/GoogleScannedObjects>. 6
- [54] Gernot Riegler and Vladlen Koltun. Free view synthesis. *ECCV*, 2020. 2
- [55] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. *CVPR*, 2019. 2
- [56] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. *CVPR*, 2016. 2, 6
- [57] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *NeurIPS*, 33, 2020. 2
- [58] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *NeurIPS*, 2020. 2
- [59] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. *CVPR*, 2019. 1, 2, 6, 8
- [60] Vincent Sitzmann, Michael Zollhofer, and Gordon Wetzstein. Scene representation networks: Continuous 3D-structure-aware neural scene representations. *NeurIPS*, 2019. 2, 6, 8
- [61] 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. *arXiv*, 2020. 2
- [62] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. *CVPR*, 2019. 2
- [63] Tiancheng Sun, Zexiang Xu, Xiuming Zhang, Sean Fanello, Christoph Rhemann, Paul Debevec, Yun-Ta Tsai, Jonathan T. Barron, and Ravi Ramamoorthi. Light stage super-resolution: Continuous high-frequency relighting. *SIGGRAPH Asia*, 2020. 4
- [64] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 2, 7
- [65] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM TOG*, 2019. 2
- [66] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d scene representation and rendering. *arXiv preprint arXiv:2010.04595*, 2020. 2
- [67] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. *CVPR*, 2017. 2
- [68] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 12
- [69] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017. 4, 12

- [70] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. Disn: Deep implicit surface network for high-quality single-view 3d reconstruction. *NeurIPS*, 2019. 2
- [71] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction with implicit lighting and material. *NeurIPS*, 2020. 2
- [72] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. *arXiv preprint arXiv:2012.02190*, 2020. 2
- [73] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. *CVPR*, 2018. 6
- [74] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *SIGGRAPH*, 2018. 1, 2, 6

6. Appendix

A. Additional implementation details

Feature extraction network architecture. To render a target view, our system takes a set of source views and extracts their features using a network with shared weights. We implement the feature extraction network using a U-Net-like architecture, where the encoder is adapted from ResNet34 [16] as implemented in PyTorch [46]. We replace all Batch Normalization [22] with Instance Normalization [68] as in [37], and remove max-pooling and use instead strided convolutions. Our network is fully convolutional and accepts input images of variable size. We take a single image of size $640 \times 480 \times 3$ as an example input and present a detailed network architecture in Tab. 4. Our code and model will be made available.

Input (id: dimension)	Layer	Output (id: dimension)
0: $640 \times 480 \times 3$	7×7 Conv, 64, stride 2	1: $320 \times 240 \times 64$
1: $320 \times 240 \times 64$	Residual Block 1	2: $160 \times 120 \times 64$
2: $160 \times 120 \times 64$	Residual Block 2	3: $80 \times 60 \times 128$
3: $80 \times 60 \times 128$	Residual Block 3	4: $40 \times 30 \times 256$
5: $40 \times 30 \times 256$	3×3 Upconv, 128, factor 2	6: $80 \times 60 \times 128$
[3, 6]: $80 \times 60 \times 256$	3×3 Conv, 128	7: $80 \times 60 \times 128$
7: $80 \times 60 \times 128$	3×3 Upconv, 64, factor 2	8: $160 \times 120 \times 64$
[2, 8]: $160 \times 120 \times 128$	3×3 Conv, 64	9: $160 \times 120 \times 64$
9: $160 \times 120 \times 64$	1×1 Conv, 64	Out: $160 \times 120 \times 64$

Table 4: **Feature extraction network architecture.** ‘Conv’ stands for a sequence of operations: convolution, rectified linear units (ReLU) and Instance Normalization [68]. ‘Upconv’ stands for a bilinear upsampling with certain factor, followed by a ‘Conv’ operation with stride 1. ‘[·, ·]’ represents channel-wise concatenation of two feature maps. The residual blocks have similar structures to those in the original ResNet34 [16] design, except that the first residual block has stride equal to 2 and all Batch Normalization layers are replaced with Instance Normalization layers. The output 64-dimensional feature map will be split into two feature maps of 32 dimension, which are then used as inputs to the coarse and fine IBRNet, respectively.

IBRNet network architecture. Fig. 5 shows the detailed network architecture of IBRNet and the process of predicting the volume density and color at a single 5D location. To obtain the density at each 5D location, we first aggregate the multi-view features drawn from source views to obtain a density feature that encodes the density information for this point. We then aggregate the density information of all samples on the ray to enable visibility reasoning for better density prediction. We implement this ray-wise operation using our proposed module called ‘ray transformer’. Specifically, we first apply positional encoding [69] to all density features on the ray, so that the network is aware of the spatial ordering for samples on the ray. We then use a single multi-head self-attention [69] layer to incorporate long-range contextual information. Separately, to obtain the

view-dependent color, we predict a set of blending weights to blend the image colors drawn from source views. IBRNet is invariant to permutations of source views, and supports a variable number of source views as well as samples on a ray.

Training details. At training time, we sample points in space and project them in the source views to fetch the corresponding colors and image features. However, the projected pixel for a sample may be located outside of the image plane. In this case, we discount this source view for this sample. If a point is not projected into the image plane of any source view, we set the volume density of this point to be zero. If there are less than 3 samples on a ray that have valid density values, we ignore this ray in the loss function during training.

For pre-training, We train on eight V100 GPUs with a batch size of 3,200 to 9,600 rays depending on image resolution and the number of source views. Within each batch, we sample rays from eight different scenes. Within each scene, we sample rays randomly from a single image for training. Pre-training takes about a day to finish. For fine-tuning, the time needed is scene-dependent. It takes us about 6 hours on a single V100 GPU to achieve the reported performance for each of the *Real Forward-Facing* [41] scenes. The fine-tuning stage may be accelerated with better hyperparameters for optimization.

B. Additional qualitative results

Ray transformer. We provide qualitative comparison of our model trained with and without ray transformer in Fig. 6. Fig. 6 shows that our proposed ray transformer significantly improves the synthesis quality especially for challenging regions (e.g., near occlusion boundaries).

Geometry Visualization. We visualize the proxy geometry by accumulating the predicted density values on each ray. Fig. 7 shows that our pretrained model produces reasonable proxy geometry, and the quality of the synthesized images and the underlying geometry improves when our model is fine-tuned.

Qualitative results on *Realistic Synthetic* 360° [42]. We provide qualitative results of our pretrained and finetuned model on *Realistic Synthetic* 360° [42] in Fig. 8. The last column in Fig. 8 shows a challenging scenario where our method does not work well.

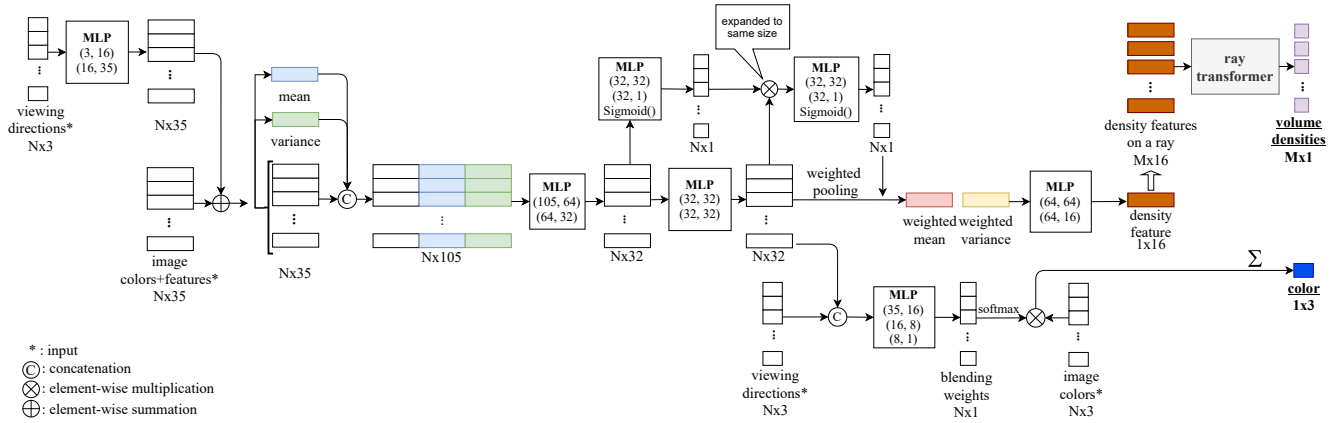


Figure 5: **IBRNet architecture**. N denotes the number of source views and M is the number of samples on a ray. For each “MLP” box, (i, j) represents a linear layer with input dimension i and output dimension j . ELU is used between each two adjacent linear layers as the non-linear activation function. When “MLP” takes in a stack of feature vectors (i.e., feature vectors from all source views), the MLP is applied to each feature vector with shared weights. “weighted pooling” computes the weighted mean and variance of the $N \times 32$ feature vector using the learned $N \times 1$ weight vector. The ray transformer module contains a single multi-head self-attention layer with the number of heads set to 4. Viewing directions shown in the figures are relative viewing directions, i.e., the viewing direction of the query ray relative to the viewing directions from source views.



Figure 6: **Qualitative comparison of our model trained with and without ray transformer**. The first and second columns show the results of our pretrained model without and with the ray transformer module, respectively. The last column shows the ground truth images. Without ray transformer, the synthesized images exhibit severe “black hole” artifacts especially near occlusion boundaries where the network fails to infer surface locations correctly. Ray transformer eliminates such artifacts by enabling more informed density prediction.

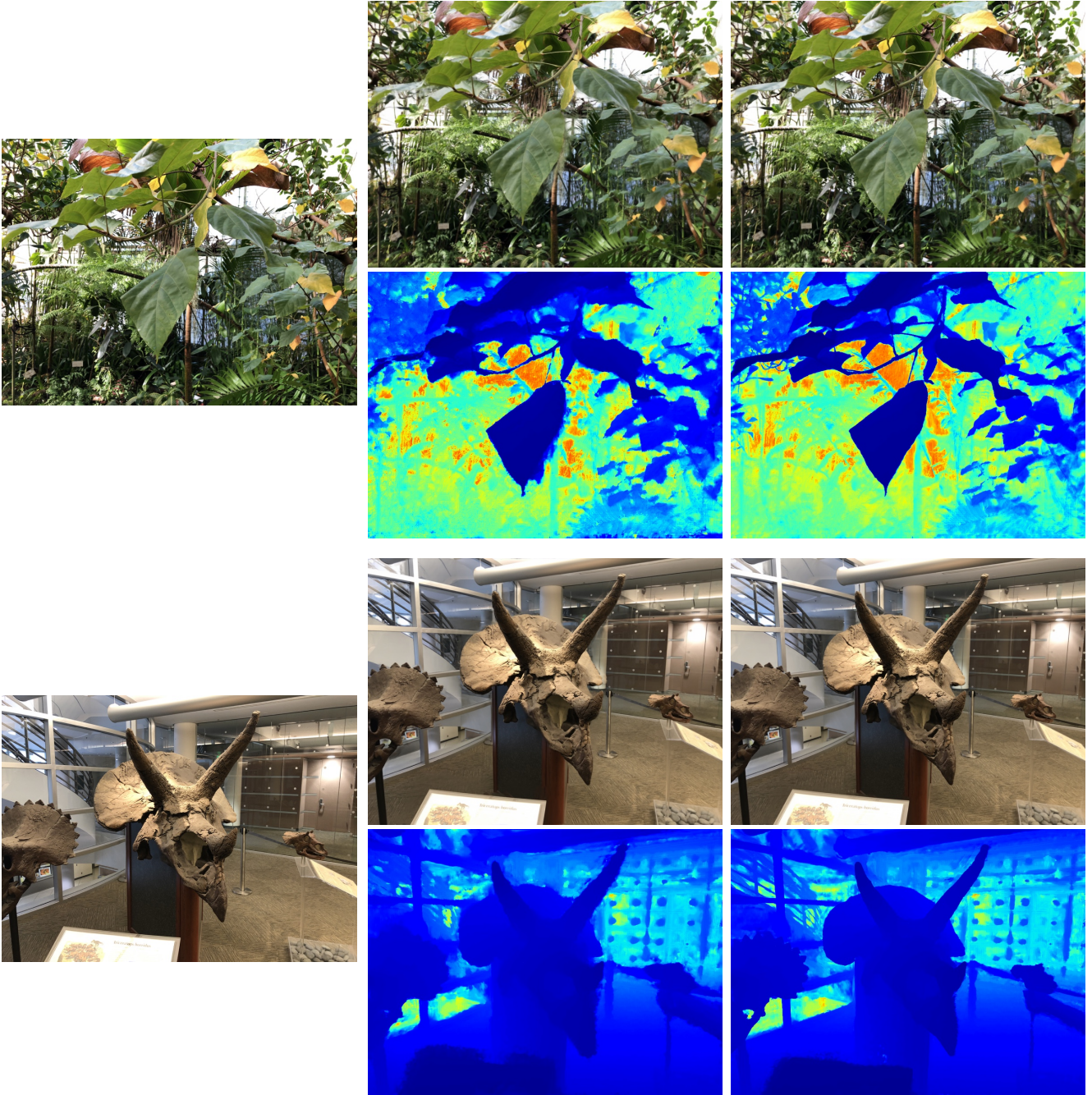


Figure 7: **Geometry Visualization.** We visualize the proxy geometry and synthesized images generated by our pretrained and fine-tuned models for two scenes *leaves* and *horns*. For each scene, the first column shows the ground truth image. The second column shows the results, i.e., synthesized image (top) and depth map (bottom), using our pretrained model. the last column shows the results of our model after finetuned on each scene.

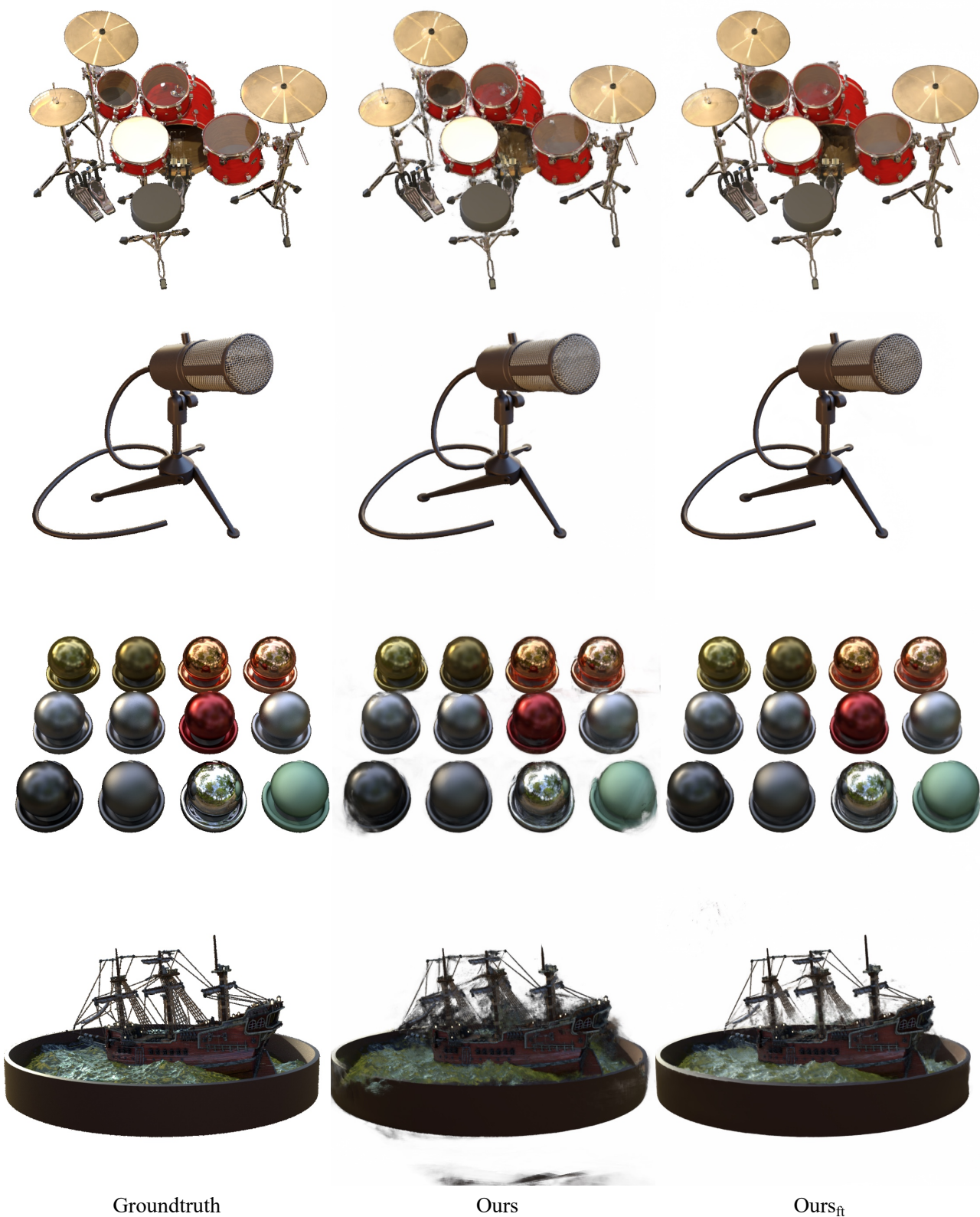


Figure 8: **Qualitative results on *Realistic Synthetic* 360° [42].** The first column shows the ground truth images. The second column shows the synthesized images without per-scene fine-tuning. The last column shows the synthesized images with per-scene fine-tuning. The last row is a failure case due to sparse source views and complex geometry.