

# Neural Pixel Composition for 3D-4D View Synthesis from Multi-Views

Aayush Bansal    Michael Zollhoefer  
 Reality Labs Research, Pittsburgh, USA

<https://www.aayushbansal.xyz/npc>



Figure 1. Our approach allows us to capture small details better than existing methods. We show novel views (top-row) synthesized using our approach and zoom on the details for each view (bottom-row). Our model is trained for 10 minutes. **Best viewed in electronic format.**

## Abstract

We present *Neural Pixel Composition (NPC)*, a novel approach for continuous 3D-4D view synthesis given a discrete set of multi-view observations as input. Existing state-of-the-art approaches require dense multi-view supervision and an extensive computational budget. The proposed formulation reliably operates on sparse and wide-baseline multi-view images/videos and can be trained efficiently within a few seconds to 10 minutes for hi-res (12MP) content. Crucial to our approach are two core novelties: 1) a representation of a pixel that contains color and depth information accumulated from multi-views for a particular location and time along a line of sight, and 2) a multi-layer perceptron (MLP) that enables the composition of this rich information provided for a pixel location to obtain the final color output. We experiment with a large variety of multi-view sequences, compare to existing approaches, and achieve better results in diverse and challenging settings.

## 1. Introduction

Novel views can be readily generated if we have access to the underlying 6D plenoptic function  $R(\theta, \mathbf{d}, \tau)$  [1, 23] of the scene that models the radiance incident from direction  $\theta \in \mathbb{R}^2$  to a camera placed at position  $\mathbf{d} \in \mathbb{R}^3$  at time  $\tau$ . Currently, no approach exists that can automatically reconstruct an efficient space- and-time representation of the plenoptic function given only a (potentially sparse) set of multi-view measurements of the scene as input. The core idea of image-based rendering [22, 38] is to generate novel views based on re-projected information from a set of calibrated source views. This re-projection requires a high-quality estimate of the scene’s geometry and is only correct for Lambertian materials, since the appearance of specular surfaces is highly view-dependent. Building a dense 3D volume from multi-view inputs that provides correct 3D information for a pixel is a non-trivial task.

Recent approaches such as Neural Radiance Fields (NeRF) [27] and Neural Volumes (NV) [20] attempt to create rich 3D information along a ray of light by sampling 3D

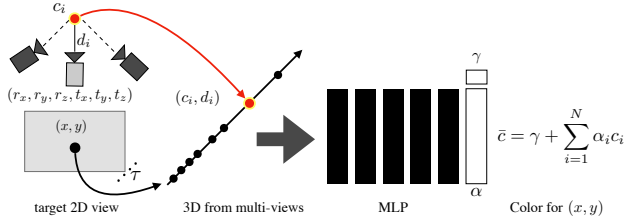


Figure 2. **Color for a pixel location:** Our goal is to estimate the color for every pixel location  $(x, y)$  for a time  $\tau$  given camera extrinsic parameters  $(r_x, r_y, r_z, t_x, t_y, t_z)$ . We collect a rich 3D descriptor consisting of color ( $c$ ) and depth ( $d$ ) information from multiple stereo-pairs using an off-the-shelf disparity estimation module [47]. We learn a multi-layer perceptron (MLP) to compose color and depth. The final output color  $\bar{c}$  is obtained by a simple dot-product of a blending weight  $\alpha$  (output of MLP) and the corresponding color samples.  $\gamma$  is a regressed color correction term per pixel.

points at regular intervals given a min-max bound. Radiance fields are highly flexible 3D scene representations that enables them to represent a large variety of scenes including semi-transparent objects. The price to be paid for this flexibility is that current approaches are restricted to datasets that provide dense 3D observations [20, 27, 32–34, 49], can only model bounded scenes [5, 20, 25, 27, 44, 48], and require intensive computational resources [20, 27, 49]. In contrast, we introduce a multi-view composition approach that combines the insights from image-based rendering [39] with the power of neural rendering [42] by learning how to best aggregate information from different views given only imperfect depth estimates as input. Figure 1 shows novel views synthesized using our approach for different multi-view sequences and the reconstructed details for each example.

We accumulate rich 3D information (color and depth) for a pixel location using an off-the-shelf disparity estimation approach [47] given multiple stereo pairs as input. We then learn a small multi-layer perceptron (MLP) for a given multi-view sequence that inputs the per-pixel information at a given camera position and outputs color at the location. Figure 2 illustrates the components of our approach. We train an MLP for a sequence by sampling random pixels given multi-views. In our experiments, we observe that a simple 5-layer perceptron is sufficient to generate high-quality results. Our model roughly requires 1 GB of GPU memory and can be trained within a few seconds to 10 minutes from scratch. The trained model allows us to perform a single forward-pass at test time for each pixel location in a target camera view. A single forward pass per pixel is more efficient than radiance field based approaches that require hundreds of samples along each ray. Finally, the alpha values ( $\alpha_i$ ) allow us to perform dense 3D reconstruction of the scene by selecting appropriate depth values at a given

location. In summary, our **contributions** are:

- A surprisingly simple, yet effective approach that requires limited computational resources for novel view synthesis from calibrated multi-view images or videos. The proposed method allows us to synthesize novel views given sparse unconstrained multi-views, where existing state-of-the-art approaches struggle.
- Our approach offers a natural extension to the 4D view synthesis problem. Our approach is also able to obtain dense depth map and 3D reconstruction on challenging in-the-wild scenes.
- Our approach enables us to reconstruct small details in the scene better than existing methods. Finally, we study the generalizability of the learned model using hi-res studio captures. We observe that the model learned on a single time-instant for one subject generalizes to unseen time instances and unseen subjects without any fine-tuning.

## 2. Related Work

Our novel view synthesis work is closely related to several research domains, such as classical 3D reconstruction and plenoptic modeling, as well as neural rendering for static and dynamic scenes. In the following, we cover the most related approaches. For a detailed discussion of neural rendering approaches, we refer to the surveys [39, 42, 43].

**Plenoptic Modeling and NeRF:** Plenoptic function [1, 23] does not require geometric modeling. A plenoptic or a light-field camera [10, 18, 28] captures all possible rays of light (in a bounded scene), which in turns enables the synthesis of a new view via a per-ray look-up. Recent approaches such as NeRF [27] and follow-up work [4, 46, 49] employ a multi-layer perceptron (MLP) that infers color and opacity values at 3D locations along each camera ray. These color and opacity values along the ray are then being integrated to obtain the final pixel color. This requires: 1) dense multi-view inputs [5, 48]; 2) perfect camera parameters [14, 19]; and 3) a min-max bound to sample 3D points along a ray of light [33, 49]. We observe degenerate outputs if all three conditions are not met (as shown in Figure 3). Different approaches either use prior knowledge or a large number of multi-view sequences [5, 35, 44, 48], additional geometric optimization [14, 19, 29], or large capacity models to separately capture foreground and background [49]. We use an off-the-shelf disparity estimation module [47] that allows us to accumulate 3D information for a given pixel. A simple MLP provides us with blending parameters that enable the composition of color information. This allows us to overcome the above-mentioned three challenges albeit using limited computational resources to train/test the model.

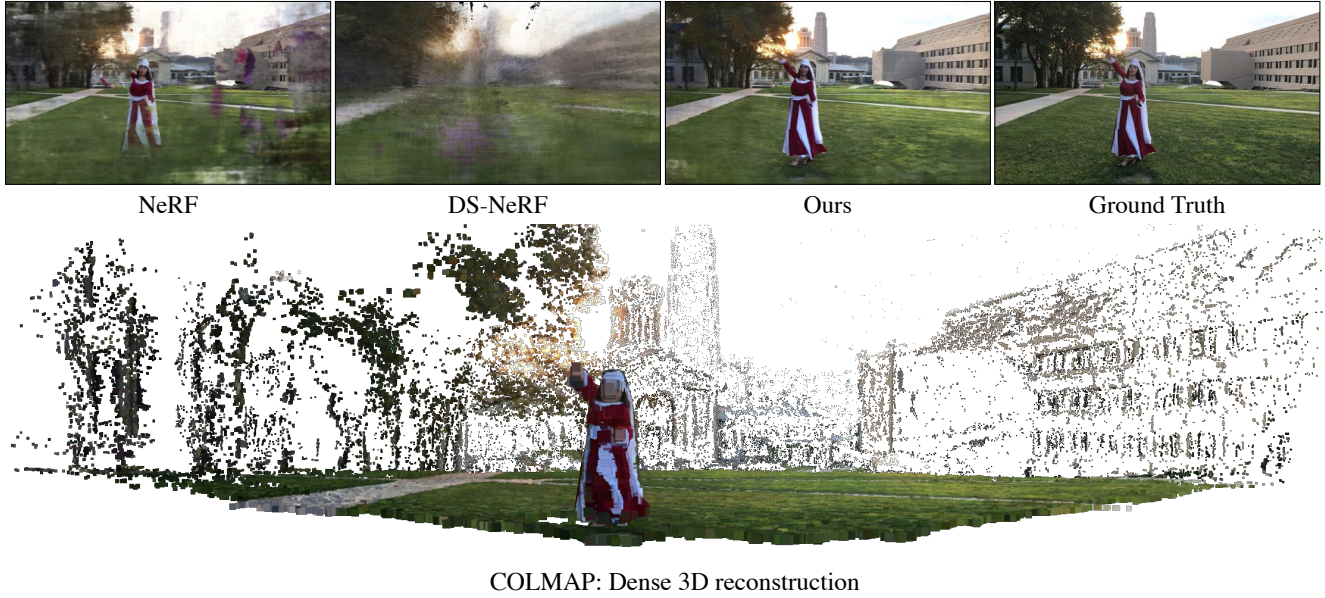


Figure 3. **View synthesis given sparse and spread-out multi-views:** Our approach allows us to operate on sparse multi-views of unbounded scenes [3]. We show novel view points for a fixed time instant for three unbounded scenes. Prior approaches such as NeRF [27] and DS-NeRF [6] lead to degenerate outputs on these sequences. We also show the 3D reconstruction using COLMAP [36, 37] for the sequence in the top-row. We observe that dense 3D reconstruction from sparse views is non-trivial for COLMAP.

**3D Reconstruction and View Synthesis:** Another approach to solve the problem is to obtain dense 3D reconstruction from the input images [11] and project 3D points to the target view. There has been immense progress in densely reconstructing the static world from multi-view imagery [8, 15], internet scale photos [2, 12, 36, 41], and videos [37]. Synthesizing a novel view from accumulated 3D point clouds may not be consistent due to varying illumination, specular material, and different cameras used for the capture of the various viewpoints. Riegler et al. [33, 34] use a neural network to obtain consistent visuals given a dense 3D reconstruction. This works well for dense multi-view observations [17]. However, 3D reconstruction is sparse given wide-baseline views or scenes with specular surfaces. This is highlighted in Figure 3, which shows 3D reconstruction results of COLMAP [36, 37] using one of the sequences. Recently, DS-NeRF [6] use sparse 3D points from COLMAP along with NeRF to learn better and faster view synthesis. As shown in Figure 3, adding explicit depth information enables DS-NeRF to capture scene structure but still struggles with details.

**Layered Depth and Multi-Plane Images:** Closely related to our work are layered depth images [24, 26, 30, 31, 38, 51] that learn an alpha composition for multi-plane images at discrete depth positions. In this work, we did not restrict our approach to 2D planes or specific depth locations. Instead, we learn a representation for a pixel at arbitrary depth locations. A pixel-wise representation not only allows us

to interpolate, but also to extrapolate, and obtain dense 3D reconstruction results. Since we have a pixel-wise representation, we are able to generate 12MP resolution images without any modifications of our approach. Prior work has demonstrated results on a max 2MP resolution content.

**4D View Synthesis:** Most approaches are restricted to 3D view synthesis [26, 27] and would require drastic modifications [7, 32] to be applied to the 4D view-synthesis problem. Lombardi et al. [21] employ a mixture of animated volumetric primitives to model the dynamic appearance of human heads from dense multi-view observations. Open4D [3] requires foreground and background modeling for 4D visualization. Our work does not require major modifications to extend to 4D view-synthesis. In addition, we do not require explicit foreground-background modeling for 4D view synthesis. We demonstrate our approach on the challenging Open4D dataset [3] where the minimum distance between two cameras is 50cm. Our composition model trained on a single time instant also enables us to do 4D visualization for unseen time instances. Finally, the model learned for view synthesis enable dense depth map and 3D reconstruction from multi-views.

### 3. Method

We are given  $M$  multi-view images with camera parameters (intrinsics and extrinsics) as input. Our goal is to learn a function,  $f$ , that inputs pixel information ( $\mathbf{p}$ ),  $\mathbf{p} \in \mathbb{R}^{N_p}$ , and outputs color ( $\bar{\mathbf{c}} \in \mathbb{R}^3$ ) at that location, i.e.,  $f : \mathbf{p} \rightarrow \bar{\mathbf{c}}$ .

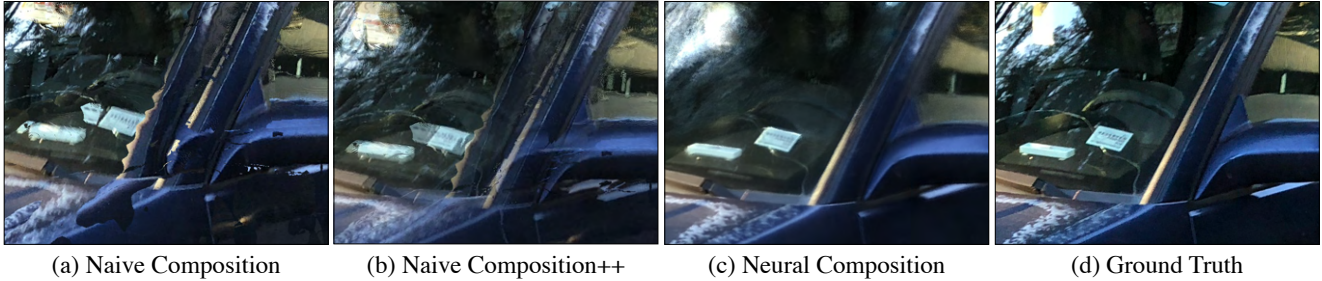


Figure 4. **Naive Composition vs. Neural Composition:** The baseline naively uses multiple stereo-pairs to generate the final output (**Naive Composition**). (a) For each pixel location, we select the color value for the closest depth location. (b) We also take the average of color values for the three closest depth locations (**Naive Composition++**). (c) We contrast these results with **Neural Composition** which uses an MLP to compose the color values. We observe that the MLP nicely composes the color values despite noisy depth estimates and fills the missing regions. (d) We also show ground truth for reference. [Best viewed in electronic format.](#)

Learning such a function is challenging since we live in a 3D-4D world and images provide only 2D measurements. We present two crucial components: 1) a representation of a **pixel** that contains relevant multi-view information for high-fidelity view synthesis; and 2) a multi-layer perceptron (**MLP**) that inputs the pixel information and outputs the color.

**Overview:** We input a pixel location  $(x, y)$  given corresponding camera parameters  $(r_x, r_y, r_z, t_x, t_y, t_z)$  at time,  $\tau$ , along with an array of possible 3D points along the line of sight. The  $i^{\text{th}}$  location of this array contains depth ( $d_i$ ) and color ( $c_i$ ). The MLP outputs alpha ( $\alpha_i$ ) values for the  $i^{\text{th}}$  location that allow us to obtain the final color at  $(x, y)$ . The MLP also outputs gamma,  $\gamma \in \mathbb{R}^3$ , which is a correction term learned by the model. We get the final color at pixel location  $(x, y)$  as:  $\bar{c} = \gamma + \sum_{i=1}^N \alpha_i c_i$ , where  $N$  is the number of points in the array. We describe our representation of a pixel in Sec. 3.1 and the MLP in Sec. 3.2.

### 3.1. Representation of a Pixel

Given a pixel location  $(x, y)$  for a camera position  $(r_x, r_y, r_z, t_x, t_y, t_z)$ , our goal is to collect dense 3D information that contains depth and color information at all possible 3D points along a line of sight. We obtain 3D points via two-view geometry [11] by forming  $\binom{M}{2}$  stereo-pairs. The estimated disparity between a stereo pair provides the depth for the 3D point locations. Multiple stereo pairs allow us to densely populate 3D points along the rays.

**Color and Depth Array:** We use multiple stereo pairs to build an array of depth (**d**) and color (**c**) for a pixel. We store the values in order of increasing depth, i.e.,  $d_{i+1} \geq d_i$ . The array is similar to a ray of light that travels in a particular direction connecting the 3D points. We limit the number of 3D points to be  $N$ . If there are less than  $N$  depth observations, we set  $d_i = 0$  and  $c_i = (0, 0, 0)$ . If there are more than  $N$  observation, we take closest  $N$  3D points.

**Uncertainty Array:** In this work, we use an off-the-shelf disparity estimation module from Yang et al. [47]. This approach provides an estimate of uncertainty (entropy) for each prediction. We also keep an array of uncertainty values ( $\mathfrak{H}$ ) of equal size as the depth array (obtained from disparity and camera parameters), s.t.,  $\mathfrak{H}_i \in [0, 1]$ , where a higher value represents higher uncertainty. The uncertainty allow us to suppress noise or uncertain 3D points.

**Encoding Spatial Information:** For each pixel, we concatenate its spatial location, i.e.,  $(x, y)$  location and camera position  $(r_x, r_y, r_z, t_x, t_y, t_z)$ . We employ high-frequency positional encoding [26] to represent spatial information of a pixel for a given camera position. We normalize the pixel coordinates, s.t.,  $x \in [-1, 1]$  and  $y \in [-1, 1]$ .

**Incorporating Temporal Information:** Our approach enables a natural extension to incorporate temporal information. Given a temporal sequence with  $T$  frames, we represent each time instant as a Gaussian distribution with peak at the frame  $\tau$ . We concatenate the color, depth, and uncertainty array alongside the spatial and temporal information in a single  $N_p$ -dimensional array, where  $N_p$  is sum of the dimensions of each term. We input this array to the MLP to compute the color output at the pixel location.

### 3.2. Multi-Layer Perceptron (MLP)

Our goal is to output blending values  $\alpha$  that enable us to take the appropriate linear combination of color values in the color-array. A naive way is to directly use the output of the last layer of the MLP as an alpha array and compute a dot product with the color-array:

$$f(x, y, \tau, r_x, r_y, r_z, t_x, t_y, t_z, \mathbf{c}, \mathbf{d}, \mathfrak{H}) = \mathbf{w}. \quad (1)$$

While this is reasonable, it assumes that the MLP will implicitly understand the relationship between color (**c**), depth

( $d$ ), and uncertainty ( $\mathfrak{H}$ ). This is challenging to learn. In this work, we observe that explicitly using the depth and uncertainty with the output of the MLP ( $w$ ) enables better view synthesis. We, therefore, define:

$$\alpha_i = \frac{(1 - \mathfrak{H}_i)e^{-(w_i d_i - \mu)^2}}{\sum_{j=1}^N (1 - \mathfrak{H}_j)e^{-(w_j d_j - \mu)^2}}, \quad (2)$$

where  $\mu = \frac{1}{N} \sum_{j=1}^N w_j d_j$ . The Gaussian distribution forces the model to select color values belonging to depth location that are: 1) closest to the average depth value; and 2) are confident and less noisy. We employ these alpha values together with the original color array to predict the final values ( $\bar{c}$ ):

$$\bar{c} = \sum_{i=1}^N \alpha_i c_i + \gamma, \quad (3)$$

where  $\gamma$  is an additional correction term that helps us to obtain sharp outputs. Note that the fifth layer of the MLP outputs  $\alpha$  and  $\gamma$  values.

**Multi-Layer Perceptron:** We employ a 5-layer perceptron. Each linear function has 256 activations followed by a non-linear ReLU activation function. We train the MLP in a self-supervised manner using a photometric  $\ell_1$ -loss:

$$\min_f \mathcal{L} = \sum_{k=1}^m \|c_k - \bar{c}_k\|_1, \quad (4)$$

where  $c_k$  and  $\bar{c}_k$  are the ground truth color and predicted color respectively for the  $k^{th}$  pixel, and  $m$  is the number of randomly sampled pixels from the  $M$  images. We train the MLP from scratch using the Adam optimizer [16]. We randomly sample 4 images, and sample 256 pixels from each image for every forward/backward pass. The learning rate is kept constant at 0.0002 for the first 5 epochs and is then linearly decayed to zero over next 5 epochs. We observe that composition model converges around in a few seconds of training on a single GPU with 1 GB GPU memory.

**Naive Composition:** One can also naively use the pixel representation to generate the final output by selecting the color value for the closest depth location. A slightly nuanced version is to take average of color values for three closest depth location (**Naive Composition++**). We use this naive composition for comparisons in our work. Figure 4 shows the importance of using neural composition via MLP over naive composition. We believe it is an importance baseline for view synthesis as this simple nearest-neighbor based method generates results without any training.

## 4. 3D View Synthesis

We study various aspects of 3D view synthesis using our approach: (1) synthesizing novel views given sparse and

24 sequences	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS [50] $\downarrow$
LLFF [26]	15.187 $\pm$ 2.166	0.384 $\pm$ 0.082	0.602 $\pm$ 0.090
NeRF [27]	13.693 $\pm$ 2.050	0.317 $\pm$ 0.094	0.713 $\pm$ 0.089
DS-NeRF [6]	14.531 $\pm$ 2.603	0.316 $\pm$ 0.099	0.757 $\pm$ 0.040
DS-NeRF** [6]	15.346 $\pm$ 2.276	0.389 $\pm$ 0.076	0.716 $\pm$ 0.048
Naive Composition	15.480 $\pm$ 1.928	0.372 $\pm$ 0.061	0.665 $\pm$ 0.065
Naive Composition++	16.244 $\pm$ 2.186	0.442 $\pm$ 0.074	0.616 $\pm$ 0.063
Ours	<b>17.946 <math>\pm</math> 1.471</b>	<b>0.562 <math>\pm</math> 0.077</b>	<b>0.534 <math>\pm</math> 0.061</b>

Table 1. **Sparse and Unconstrained Multi-Views:** We evaluate on the 24 sparse and unconstrained multi-view sequences of the Open4D dataset [3]. We train NeRF [27] and DS-NeRF [6] models for each sequence. DS-NeRF [6] employs additional depth along with NeRF. We trained two versions of DS-NeRF. One where we use the same model as NeRF with additional depth supervision. The second version is DS-NeRF\*\* with tuned hyperparameters. We also use off-the-shelf LLFF model. We observe degenerate outputs using LLFF, NeRF and DS-NeRF, especially for unbounded scenes. However, our approach is able to reliably generate novel views in twenty times less time. Training a NeRF or a DS-NeRF model takes roughly 420 minutes per sequence whereas our approach take 10 minutes (including pre-processing multi-view content). We also generate results using naive composition and obtain better results than prior work.

unconstrained multi-views (Sec. 4.1); (2) synthesizing hires 12MP content (Sec. 4.2); and (3) scenes with unbounded depth and influence of the number of views (Sec. 4.3). We demonstrate our approach on hi-resolution studio capture in Sec. A.4 and observe that our approach can generalize to unseen subjects and unknown time instances from a single time-instant. We study convergence in Sec. A.5 where we observe that our model gets close to convergence within a few seconds of learning.

### 4.1. Sparse and Unconstrained Multi-Views

We use 24 sequences of sparse and unconstrained real-world samples from the Open4D dataset [3]. Open4D consists of temporal sequences. We use certain time instants for 3D view synthesis. The minimum distance between two adjacent cameras is 50cm in these sequences. We contrast our approach with NeRF [27]. We also study DS-NeRF [6], which additionally employs sparse 3D point clouds from COLMAP for training NeRF. DS-NeRF has shown promising results given the sparse views from LLFF dataset [26]. Both approaches are trained for 200, 000 iterations (roughly 420 minutes) per sequence on a NVIDIA V100 GPU. Table 1 compares the performance of different methods on held-out views from these sequences using PSNR, SSIM, and LPIPS (AlexNet) [50]. In this work, we observe that these three evaluation criteria are not self-sufficient in determining the relative ranking of different methods. While PSNR and SSIM may favor smooth or blurry results [50], LPIPS may ignore the structural consistency in images. We

8 sequences	PSNR $\uparrow$	MCSSIM $\uparrow$	LPIPS $\downarrow$
NeRF [27]	22.009 $\pm$ 3.148	0.757 $\pm$ 0.156	0.487 $\pm$ 0.180
NeX [45]	22.292 $\pm$ 3.137	0.774 $\pm$ 0.152	0.423 $\pm$ 0.156
Naive Composition	16.624 $\pm$ 2.906	0.648 $\pm$ 0.197	0.342 $\pm$ 0.096
Naive Composition++	17.535 $\pm$ 2.698	0.688 $\pm$ 0.184	0.317 $\pm$ 0.107
<b>Ours (10 minutes)</b>	<b>22.868 <math>\pm</math> 4.588</b>	<b>0.802 <math>\pm</math> 0.140</b>	<b>0.269 <math>\pm</math> 0.120</b>

Table 2. **Shiny dataset:** We study our approach on the 8 real sequences from the Shiny dataset [45]. We use NeRF [27] and NeX [45] results from the authors of Shiny dataset [45]. NeX use multiple GPUs and a few days to train a model. Our approach gets better performance in only a few minutes.

posit that it is important to look at all three criteria and not one. Figure 3 shows the qualitative performance of our approach on these challenging sequences. We observe degenerate results using NeRF on these sequences. DS-NeRF also results in degenerate outputs most of the times except for the scenes with bounded depth. Our approach is able to generate high-quality results (with details such as faces, hair, dress, etc.) in this setting both qualitatively and quantitatively. Total time taken to process (pre-processing multi-view content and training the composition model) a sequence is less than 10 minutes. We also observe that a naive pixel composition can also yield meaningful results better than prior work. However, we obtain better pixel composition using MLPs. The details of these sequences are available in Appendix A.1.

## 4.2. Better Details in Lesser Time

We use 8 multi-view sequences from the Shiny Dataset [45] that consists of multi-views captured for specular surfaces. Table 2 shows the performance of different methods. We use the results generated by NeX [45]. It takes multiple GPUs for a few days to get these results. Our model trained for 10 minutes achieves better performance. Figure 6 contrasts our method with NeX. We observe small details are better captured by our method.

We use 12 high-resolution ( $4032 \times 3024$ ) multi-view sequences from the LLFF dataset [26] that contain challenging specular surfaces. In this setting, we train NeRF [27] on these sequences for 2,000,000 iterations which take approximately 64 hours on a single NVIDIA V100 GPU (10,000 iterations take 20 minutes). Performance saturates at  $1M$  iterations after 32 hours of training. We train our model for 10 epochs, which takes around 10 minutes on a single GPU and only 1GB GPU of memory. We estimate disparity [47] for multiple stereo pairs at one-fourth resolution for these sequences. Disparity estimation using the off-the-shelf model takes less than 5 minutes per sequence on a single GPU. Table 3 contrasts the performance of NeRF models at different intervals of training using PSNR, SSIM,

12 sequences	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
NeRF (32 hours) [27]	<b>21.741 <math>\pm</math> 2.985</b>	<b>0.602 <math>\pm</math> 0.147</b>	0.584 $\pm$ 0.087
Naive Composition	16.008 $\pm$ 2.315	0.415 $\pm$ 0.142	0.427 $\pm$ 0.068
Naive Composition++	17.022 $\pm$ 2.483	0.460 $\pm$ 0.144	<b>0.406 <math>\pm</math> 0.066</b>
Ours (10 minutes)	20.953 $\pm$ 2.805	0.598 $\pm$ 0.136	0.460 $\pm$ 0.078

Table 3. **Hi-Res (12MP) View Synthesis:** We evaluate on 12 sequences from LLFF containing specular surfaces on original  $4032 \times 3024$  resolution. We contrast the performance of our approach with a NeRF model trained for  $1M$  iterations after 32 hours of training using a NVIDIA-V100 GPU. Our model converges quickly in a few minutes. Training our model requires 1 GB of GPU memory for training. The MLP allows us to obtain better results than a naive composition (Fig. 4). We refer the reader to Figure 5 for visual comparisons. We observe that details become better for NeRF when trained for long. However, our approach captures more details in a few minutes as compared to 32 hours of training of a NeRF model. Consistent with the observation of Zhang et al. [50], PSNR may favor averaged/blurry results while LPIPS favors sharp results. More analysis in Appendix A.2.

and LPIPS (AlexNet). We compute the average of per-frame statistics as the number of samples in the test set for these 12 sequences are roughly the same. We once again observe that it is crucial to include all three evaluation criteria. Figure 5 shows the results of NeRF at different intervals of time. We observe that the NeRF model improves over time and captures sharp results as suggested by LPIPS. Our method enables sharper outputs as compared to NeRF. The qualitative and quantitative analysis suggest that we can efficiently generate results on 12MP images without drastically increasing the computational resources. We also show the performance of naive composition to generate the final outputs. We observe that MLPs allow us to obtain better results. More analysis are available in Appendix A.2.

## 4.3. Unbounded Scenes and Number of Views

We study the influence of the number of views using the challenging synthetic multi-view sequences from MVS-Synth dataset [13] that consist of different unbounded scenes. We use the first 13 sequences with unbounded depth from this dataset for our analysis. Each sequence consists of 100 frames. We use 50 frames ( $1920 \times 1080$  resolution) for evaluation, and train models by varying the number of views between  $\{10, 20, 30, 40, 50\}$ . The details about train-test splits are available in Appendix A.3. For this analysis, we train 65 NeRF [27] models (each for 200,000 iterations taking roughly 420 minutes per model) and 65 models for our approach. We contrast the performance of two methods in Table 4. Without any modification, our approach can generate better results. We can also generate better results with fewer views. The camera parameters are computed using Agisoft Metashape here for each setting. We also contrast

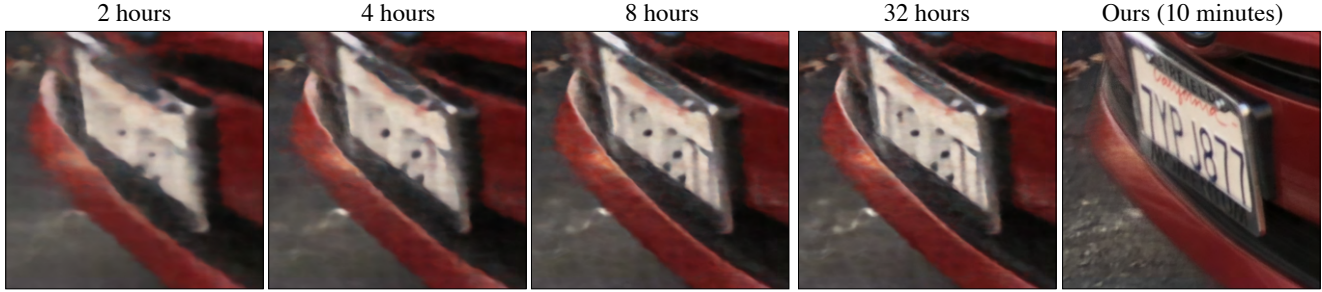


Figure 5. **Improvement in NeRF over time:** We show the progression (first 32 hours) of improvement for the NeRF model. We observe that results improve over time as details become clearer over time. Our approach can generate sharp results in only 10 minutes. One can zoom-in to see *California* written on the number plate. Full image is in the supplementary material. [Best viewed in electronic format.](#)

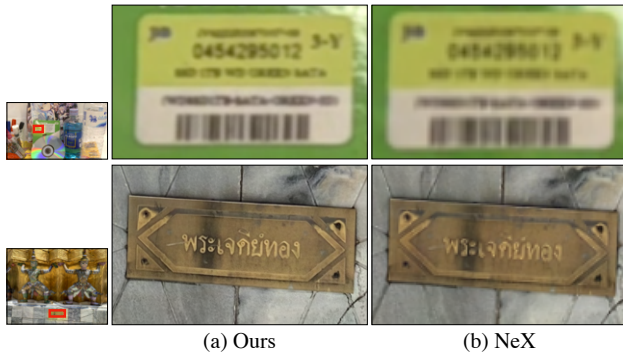


Figure 6. **Shiny Dataset:** We contrast the results of our approach (a) with (b) NeX [45] on held-out views. Our approach is able to capture the details better than NeX such as the text (0454295012 3-Y) in the top-row and the details on the plate and stone in the bottom row. [Best viewed in electronic format.](#)

13 sequences	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
num-views=10			
NeRF	16.150 $\pm$ 4.195	0.541 $\pm$ 0.139	0.619 $\pm$ 0.158
Ours	<b>18.460 <math>\pm</math> 4.099</b>	<b>0.656 <math>\pm</math> 0.129</b>	<b>0.451 <math>\pm</math> 0.167</b>
num-views=20			
NeRF	18.171 $\pm$ 4.543	0.582 $\pm$ 0.135	0.594 $\pm$ 0.171
Ours.	<b>22.414 <math>\pm</math> 4.197</b>	<b>0.766 <math>\pm</math> 0.126</b>	<b>0.289 <math>\pm</math> 0.147</b>
num-views=30			
NeRF	19.725 $\pm$ 4.759	0.619 $\pm$ 0.135	0.557 $\pm$ 0.179
Ours	<b>24.191 <math>\pm</math> 4.219</b>	<b>0.803 <math>\pm</math> 0.122</b>	<b>0.243 <math>\pm</math> 0.137</b>
num-views=40			
NeRF	20.074 $\pm$ 4.673	0.627 $\pm$ 0.132	0.556 $\pm$ 0.178
Ours.	<b>24.832 <math>\pm</math> 4.110</b>	<b>0.822 <math>\pm</math> 0.117</b>	<b>0.218 <math>\pm</math> 0.125</b>
num-views=50			
NeRF	20.244 $\pm$ 4.611	0.631 $\pm$ 0.129	0.556 $\pm$ 0.178
Ours	<b>25.529 <math>\pm</math> 4.212</b>	<b>0.836 <math>\pm</math> 0.112</b>	<b>0.198 <math>\pm</math> 0.116</b>

Table 4. **Synthetic Multi-View Sequences:** We vary the number of views in our analysis. Our approach achieves better performance even with a few views and improve as we increase the number of views.

the performance of our approach using ground truth camera parameters with estimated parameters in Appendix A.3.



Figure 7. **4D View Synthesis:** Our approach allows us to get better facial details than Open4D [3] on challenging sequences.

## 5. 4D View Synthesis

We study the ability of our approach to perform 4D view synthesis. We train our model on the temporal sequences (1920  $\times$  1080 resolution) from the Open4D dataset and contrast our approach with their method [3]. We conduct two experiments: (1) held-out temporal sequences; and (2) held-out camera views.

**Held-out temporal sequences:** We study the performance of the trained model on unseen temporal sequences. We train the model without temporal constraint. Our goal is to study the compositional ability of our model in contrast to the more explicit Open4D. The model is trained with multi-views available for 300 – 400 time instances and evaluated on unseen 100 time instances. Table 5 contrasts the performance of our approach with Open4D. Quantitatively, we observe similar performance of our approach as compared to Open4D on unseen temporal sequences. Our approach is able to capture details such as human faces consistently better than Open4D as shown in Figure 7.

**Held-out camera views:** We study the performance on unseen camera views but a known temporal sequence. We train the model for 500 time instances with and without

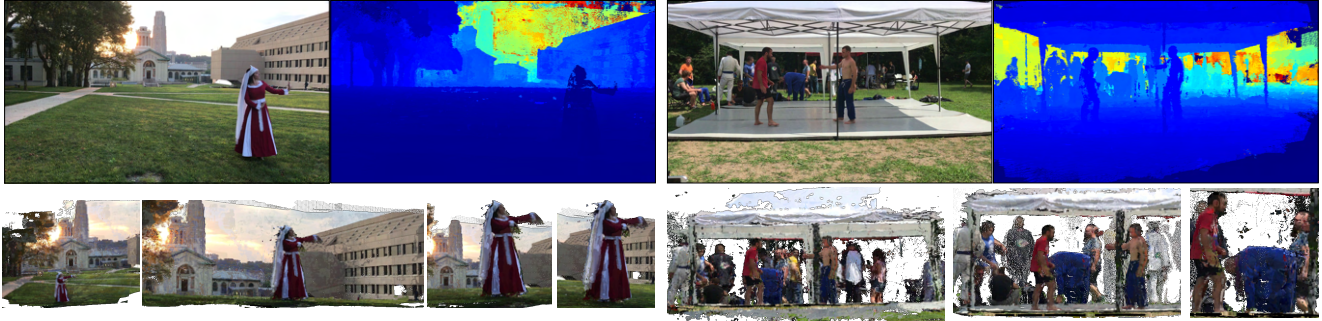


Figure 8. We show depth map (top-row) and dense 3D point clouds (bottom-row) for two sequences. The “jet blue” color corresponds to missing depth values for these images (e.g., the bottom right edge on the depth map of the first image).

5 sequences	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Naive Composition	13.723 $\pm$ 2.759	0.342 $\pm$ 0.110	0.665 $\pm$ 0.113
Open4D [3]	20.355 $\pm$ 4.425	0.626 $\pm$ 0.131	<b>0.306 <math>\pm</math> 0.079</b>
Ours	<b>21.458 <math>\pm</math> 4.690</b>	<b>0.645 <math>\pm</math> 0.145</b>	0.431 $\pm$ 0.139

Table 5. **Unseen Temporal Sequences:** We study the compositional ability of our model in contrast to the more explicit Open4D. The model is trained with multi-views available for 300 – 400 time instances and evaluated on unseen 100 time instances. There are a total of 5297 frames used for evaluation. Our approach is able to generate results competitive to Open4D without any modification.

temporal constraint to understand its importance. Table 6 contrasts the performance of our approach with Open4D. Without any heuristics and foreground-background estimation, we are able to learn a representation that allows 4D view synthesis. Our approach use a simple reconstruction loss whereas Open4D use an additional adversarial loss [9]. Using the adversarial loss enables Open4D to generate overall sharp results that leads to lower LPIPS score. However, the details are not consistent as shown in Appendix B.

## 6. Depth Map and Dense 3D reconstruction

We use the learned MLPs to construct depth map for a given view. Given an array of depth values for a pixel, we select the depth value corresponding to the maximum  $\alpha_i$  value. Multiple stereo pairs also provide us with dense 3D point clouds. However, correspondences can still be noisy, and using them with noisy camera parameters leads to poor 3D estimates. We observe that the learned MLP enables us to select good 3D points per view that can be accumulated across multi-views to obtain a dense 3D reconstruction. For each pixel, we take the top-3  $\alpha_i$  values and check if the corresponding  $d_i$  values are in the vicinity of each other (this is done by empirically selecting a distance threshold). If they are, then we select the 3D point from a stereo pair corresponding to the maximum  $\alpha_i$  value. The process is repeated for all the pixels in the available multi-views. Figure 8 shows the depth map and dense 3D reconstruction. We

5 sequences	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Naive Composition	14.584 $\pm$ 3.364	0.374 $\pm$ 0.089	0.617 $\pm$ 0.064
Open4D [3]	16.681 $\pm$ 2.718	0.498 $\pm$ 0.071	<b>0.477 <math>\pm</math> 0.061</b>
Ours (w/o T)	16.665 $\pm$ 2.365	0.519 $\pm$ 0.074	0.538 $\pm$ 0.071
Ours (w/ T)	<b>16.797 <math>\pm</math> 2.523</b>	<b>0.535 <math>\pm</math> 0.080</b>	0.522 $\pm$ 0.075

Table 6. **Held-Out Camera Views:** We contrast the performance of our approach with Open4D [3] to synthesize held-out camera views. There are a total of 2092 frames used for evaluation. We achieve similar performance. We further improve performance by incorporating temporal information as an input to the model.

do not have ground truth depth values for these sequences.

## 7. Discussion

We propose a novel approach for continuous 3D-4D view synthesis from sparse and wide-baseline multi-view observations. Leveraging a rich pixel representation that consists of color, depth, and uncertainty information leads to a high performing view-synthesis approach that generalizes well to novel views and unseen time instances. Our approach can be trained within few minutes from scratch utilizing as few as 1GB of GPU memory. In this work, we strive to provide an extensive analysis of our approach in contrast to existing methods on a wide variety of settings. Our method works well on numerous settings without incorporating any task-specific or sequence-specific knowledge. We see our approach as a first step towards more efficient and general neural rendering techniques via the explicit use of geometric information and hope that it will inspire follow-up work in this exciting field.

**Acknowledgements:** We would like to thank David Forsyth for many comments and insights that were extremely helpful in designing this work.

**Disclaimer:** This academic article may contain images and/or data from sources that are not affiliated with the article submitter. Inclusion should not be construed as approval, endorsement or sponsorship of the submitter, article or its content by any such party.



## References

- [1] Edward H Adelson, James R Bergen, et al. *The plenoptic function and the elements of early vision*, volume 2. Vision and Modeling Group, Media Laboratory, Massachusetts Institute of Technology, 1991. [1](#), [2](#)
- [2] S. Agarwal, N. Snavely, I. Simon, S. M. Seitz, and R. Szeliski. Building rome in a day. In *ICCV*, 2009. [3](#)
- [3] Aayush Bansal, Minh Vo, Yaser Sheikh, Deva Ramanan, and Srinivasa Narasimhan. 4d visualization of dynamic events from unconstrained multi-view videos. In *CVPR*, 2020. [3](#), [5](#), [7](#), [8](#), [11](#), [12](#), [19](#), [20](#), [23](#), [26](#)
- [4] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. [2](#)
- [5] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (srf): Learning view synthesis from sparse views of novel scenes. In *CVPR*, 2021. [2](#)
- [6] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *CVPR*, 2022. [3](#), [5](#), [12](#), [21](#)
- [7] Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B. Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *ICCV*, 2021. [3](#)
- [8] Yasutaka Furukawa, Carlos Hernandez, et al. *Multi-view stereo: A tutorial*. Foundation and Trends in Computer Graphics and Vision, 2015. [3](#)
- [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. [8](#)
- [10] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *SIGGRAPH*, 1996. [2](#)
- [11] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. [3](#), [4](#)
- [12] Jared Heinly. *Toward Efficient and Robust Large-Scale Structure-from-Motion Systems*. PhD thesis, The University of North Carolina at Chapel Hill, 2015. [3](#)
- [13] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [6](#), [13](#)
- [14] Yoonwoo Jeong, Seokjun Ahn, Christophehr Choy, Animesh Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *ICCV*, 2021. [2](#)
- [15] Takeo Kanade and PJ Narayanan. Historical perspectives on 4d virtualized reality. In *CVPR Workshops*, 2006. [3](#)
- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [17] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Trans. Graph.*, 2017. [3](#)
- [18] Marc Levoy and Pat Hanrahan. Light field rendering. In *Annual Conference on Computer Graphics and Interactive Techniques*. ACM, 1996. [2](#)
- [19] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. [2](#)
- [20] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 2019. [1](#), [2](#)
- [21] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Trans. Graph.*, 2021. [3](#)
- [22] L McMillan. An image-based approach to three-dimensional computer graphics. *Ph. D. Dissertation, UNC Computer Science*, 1999. [1](#)
- [23] Leonard McMillan and Gary Bishop. Plenoptic modeling: An image-based rendering system. In *SIGGRAPH*, 1995. [1](#), [2](#)
- [24] Alexandre Meyer and Fabrice Neyret. Interactive volumetric textures. In *Eurographics Workshop on Rendering Techniques*, pages 157–168. Springer, 1998. [3](#)
- [25] Marko Mihajlovic, Aayush Bansal, Michael Zollhoefer, Siyu Tang, and Shunsuke Saito. KeypointNeRF: Generalizing image-based volumetric avatars using relative spatial encoding of keypoints. In *ECCV*, 2022. [2](#)
- [26] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Trans. Graph.*, 2019. [3](#), [4](#), [5](#), [6](#), [11](#), [13](#), [17](#), [19](#), [20](#), [21](#)
- [27] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [11](#), [12](#), [13](#), [19](#), [21](#), [22](#)
- [28] Ren Ng. *Light field photography with a hand-held plenoptic camera*. PhD thesis, Stanford University, 2005. [2](#)
- [29] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022. [2](#)
- [30] Manuel M Oliveira and Gary Bishop. Image-based objects. In *Proceedings of the 1999 symposium on Interactive 3D graphics*, pages 191–198, 1999. [3](#)
- [31] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Trans. Graph.*, 2017. [3](#)
- [32] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *CVPR*, 2021. [2](#), [3](#)
- [33] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, 2020. [2](#), [3](#)
- [34] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *CVPR*, 2021. [2](#), [3](#)
- [35] Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. Dense depth pri-

- ors for neural radiance fields from sparse input views. In *CVPR*, 2022. 2
- [36] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 3
- [37] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016. 3
- [38] Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski. Layered depth images. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, 1998. 1, 3
- [39] Harry Shum and Sing Bing Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing*, 2000. 2
- [40] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Neurips*, 2019. 17
- [41] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Photo tourism: Exploring photo collections in 3d. *ACM Trans. Graph.*, 2006. 3
- [42] A. Tewari, O. Fried, J. Thies, V. Sitzmann, S. Lombardi, K. Sunkavalli, R. Martin-Brualla, T. Simon, J. Saragih, M. Nießner, R. Pandey, S. Fanello, G. Wetzstein, J.-Y. Zhu, C. Theobalt, M. Agrawala, E. Shechtman, D. B Goldman, and M. Zollhöfer. State of the Art on Neural Rendering. *Computer Graphics Forum (EG STAR 2020)*, 2020. 2
- [43] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, Tomas Simon, Christian Theobalt, Matthias Niessner, Jonathan T. Barron, Gordon Wetzstein, Michael Zollhoefer, and Vladislav Golyanik. Advances in neural rendering, 2021. 2
- [44] 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. 2
- [45] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 6, 7, 11, 13, 17, 19, 20, 22, 24
- [46] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yuxiong Wang, and David Forsyth. Diver: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. In *CVPR*, 2022. 2
- [47] Gengshan Yang, Joshua Manela, Michael Happold, and Deva Ramanan. Hierarchical deep stereo matching on high-resolution images. In *CVPR*, 2019. 2, 4, 6, 11
- [48] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. Pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [49] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [50] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5, 6, 13
- [51] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM Trans. Graph.*, 2018. 3