

Learning Generalizable Light Field Networks from Few Images

Qian Li Franck Multon Adnane Boukhayma

Inria, Univ. Rennes, CNRS, IRISA, M2S, France

Abstract. We explore a new strategy for few-shot novel view synthesis based on a neural light field representation. Given a target camera pose, an implicit neural network maps each ray to its target pixel’s color directly. The network is conditioned on local ray features generated by coarse volumetric rendering from an explicit 3D feature volume. This volume is built from the input images using a 3D ConvNet. Our method achieves competitive performances on synthetic and real MVS data with respect to state-of-the-art neural radiance field based competition, while offering a 100 times faster rendering.

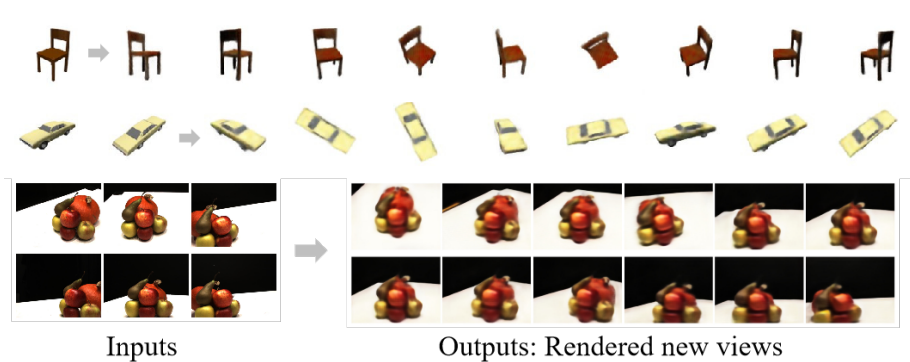


Fig. 1: Our method enables fast generation of novel views from sparse input images without 3D supervision in training.

1 Introduction

The ongoing research in computer vision and artificial intelligence has long sought to enable machines to understand and reason about the 3D world given limited observations. This ability is in fact crucial for many downstream 3D based machine learning, vision and graphics tasks. Among these, novel view synthesis is a particularly prominent problem with numerous applications in free viewpoint and virtual reality, as well as image editing and manipulation. While traditional approaches required depth information, coarse geometric proxies or dense samplings of the input views, current deep learning based approaches are relying on deep neural network’s generalization abilities across view points and 3D scenes to achieve novel view synthesis from minimal visual input.

Leveraging a training data corpus of multi-view images and camera poses, without any 3D information, our goal is to build a novel view machine that can generalize outside the training data, and extrapolate beyond the input views. Given a few input calibrated color images at test time, we expect our method to generate novel target images given new query view points (See figure 1). We are interested also in fast rendering novel view machines, that can predict new images in a single forward pass, without test time optimization. The recently popularized implicit neural representations offer numerous advantages in modelling 3D shape [1,2] and appearance [3,4] in comparison to their traditional alternatives, while being conditionable using e.g. encoders [5,6,7] and meta-learning [8,9]. In particular, Neural Radiance Fields [3] (NeRF) provide impressive novel view synthesis performances from dense input images. When coupled with convolutional encoders (e.g. [7,10]), they can additionally achieve across-scene generalization and test-time optimization free inference, in addition to reconstructing from fewer inputs. However the rendering of these methods is expensive. In fact they require sampling hundreds of 3D points along each target pixel ray, evaluating densities and view-dependent colors for all these points through a multi-layer perceptron (MLP), and building the final image through volumetric rendering of all the samples' colors and densities. Multi-scale sampling is also necessary to achieve satisfactory results.

To reduce this complexity, we propose to use an implicit neural network operating in ray space rather than the 5D Euclidean \times direction space, thus alleviating the need for per ray multi-sample evaluation and physical rendering. For a given target pixel, an MLP (i.e. light field network) maps its ray coordinate and ray features to the color directly. Key to efficient generalization, and differently from [11], we build the ray features by computing and merging 3D convolution feature volumes from the input images. These features are then rendered volumetrically into a coarse ray feature image, as illustrated in figure 2. The method is fully differentiable and trained end-to-end.

We evaluate our method on both synthetic (ShapeNet [12]) and real multi-view stereo data (DTU [13]). In the few-shot novel view optimization-free setting, we outperform comparable convolutional methods, including our own 3D convolutional baseline, and extend them to real complex data, thanks to our hybrid approach combining 3D aware ConvNets and implicit neural representations. We achieve competitive results in comparison to generalizable encoder-decoder NeRF based models, while providing orders of magnitude faster rendering (see table 2).

2 Related Work

Novel view synthesis is long-standing problem that spurred considerable work from the computer vision and graphics communities alike. While seminal work introduced various representations such as multi-plane images [14], depth layered images [15], light fields [16], depth based warping [17,18] etc., there has been a recent surge of particular interest in implicit neural shape and appearance

representations (e.g. [3,4,19,20,21,22]) and neural rendering (e.g. [23,24,25,26]). We discuss in this section existing work that we deemed most relevant to our contribution to few-shot novel view synthesis.

Early deep learning based approaches to novel view from very sparse inputs used 2D convolutional encoder-decoder architectures mapping the input to the target image, conditioned on the target view. They either predicted colors directly [27,28] or 2D flow fields [29,30,31] applied subsequently to the input. These methods were outperformed by 3D aware convolutional approaches, which rely on encoding then rendering explicit 3D latent presentations, either through volumetric rendering [32], reparameterization [33], or learnable neural rendering [34,35]. The 3D latents take the form of intrinsic scene representations [33,36,37] or extrinsic volume grids [34,35,32]. Although many of these methods could learn to generate 360-degree views from very sparse inputs especially for synthetic central object data, most of them could not scale to high resolution images, complex scenes, and real data such as multi-view stereo datasets (DTU [13]). Some also required foreground segmentation masks at training (e.g. [34]).

Implicit neural radiance fields (NeRF) [3] emerged later on as a powerful representation for novel view synthesis. It consists of an MLP mapping spatial points to volume density and view-dependant colors. Images are rendered through hierarchical volumetric rendering. It presented initially however a few limitations such as computational and time rendering complexity, requiring dense training views, lacking across-scene generalization, and requiring test-time optimization. Current research work is tackling each of these limitations (e.g. [38,39,40,7,41,42,43,44]). In particular, recent methods proposed to augment NeRFs with 2D [7,45,46] and 3D [10] convolutional features collected from the input images. Hence, they offer forward pass prediction models, i.e. test-time optimization free, while introducing generalization across scenes. However, they still need to evaluate hundreds of 3D query points per ray for inference similarly to NeRF, which makes them slow to render. Furthermore, methods such as [42,47] try to alleviate NeRFs' rendering complexity by learning view independent radiance features. In [47], the latter is combined with a single ray-dependant specular component, while Yu et al. [42] predict radiance spherical harmonic coefficients instead. Following [11], we explore here a tangent strategy, consisting in bypassing 3D implicit radiance modelling all together.

Sitzmann et al. [11] introduced recently a new implicit representation for modelling multi-view appearance. A neural light field function maps rays i.e. target pixels directly to their colors without any need for physical rendering. However, this method was not demonstrated on real MVS data (e.g. DTU [13]). It was also implemented in the auto-decoding setup, which means it requires test time optimization. The method also uses a hypernetwork for conditioning which makes it expensive to scale to bigger images in compute and memory. Differently, We propose here a more efficient local conditioning mechanism for the light field network, which allows us to scale to real data and larger images, and offers optimization-free inference.

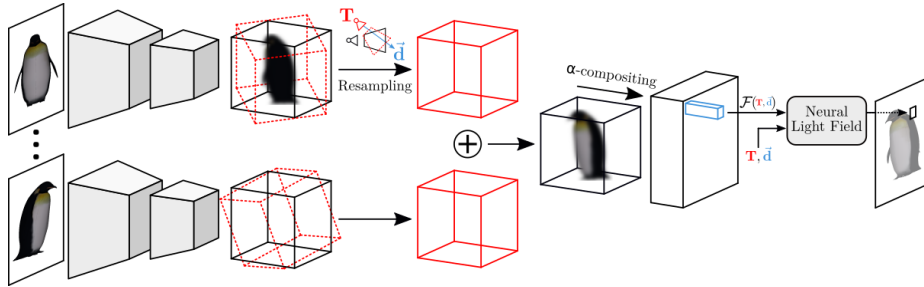


Fig. 2: Overview of our method. Given an input image, a 3D feature volume is built with a convolutional neural network (first black cube). The volume represents features inside the input view frustum. Given a target view, these features are resampled into a volume representing the target view frustum (red cube). Target feature volumes originating from different input views are aggregated using learnable weights. An image of ray features is produced by rendering the target aggregated feature volume with α -compositing. Finally the light field network maps a ray stemming from a target camera origin T and spanning a direction d , along with its convolutional feature \mathcal{F} , to the corresponding pixel color of the target image.

3 Method

Given one or few images $\{I_i\}$ of a scene or an object with their known camera parameters, i.e. camera poses $\{R_i, T_i\}$, $R_i \in SO(3)$, $T_i \in \mathbb{R}^3$, and intrinsics $K \in \mathbb{R}^{3 \times 3}$, our goal is to generate images $\{I_t\}$ for novel target views, i.e. new camera poses $\{R_t, T_t\}$. We are interested in generalization to scenes and objects unseen at training, and target views beyond input view interpolation, using merely sparse input views. We also seek fast rendering time, and we do not assume the availability of any segmentation masks neither at training nor testing.

To this end, we propose a single forward pass inference deep learning method, that uses a deep neural network to map a ray r in a projective pinhole camera model, to its desired color in the target view image c_r , using an implicit neural representation i.e. a neural light field network f . This network is conditioned with ray features \mathcal{F}_r , i.e. $c_r = f(r, \mathcal{F}_r)$. The ray features are generated through volumetric rendering of explicit 3D convolutional features built from the input images. A summary of our method is illustrated in figure 2. We present in the remaining of this section the components of the two stages of our method, namely the convolutional stage, and the neural light field network.

3.1 Feature Volume

Following seminal work (e.g. [34,35,10]), we build an explicit volume of features from an input image I_i using a fully convolutional neural network E consisting

of a succession of a 2D convolutional U-Net and several 3D convolutional blocks:

$$F_i = E(I_i), \quad (1)$$

where $I_i \in \mathbb{R}^{H \times W \times 3}$, H and W being the height and width of the input RGB image, and $F_i \in \mathbb{R}^{H_V \times W_V \times D \times C}$, H_V , W_V , D and C being respectively the height, width, depth, and the number of channels of the 3D feature volume. We refer the reader to the supplementary material for more details about the architecture of network E . The feature volume F_i is expected to encode 3D shape and appearance information of the captured object or scene in the view frustum associated with the input image, and is hence aligned pixel-wise with the latter. As we will show in the following sections, this volume will encode a prediction confidence, volume density [3], colors, in addition to more generic appearance features. We note that one limitation of these features being modelled explicitly and not implicitly as in NeRF [3] based methods is that they cannot be view direction dependent.

3.2 Feature resampling

Next, using the the input feature volume F_i aligned with the input image, we would like to create a feature volume $F_{t/i}$ aligned to the target image, that could be used subsequently to render a target feature image given the target camera pose $\{R_t, T_t\}$. Following the principles of volumetric rendering [48,3], in order to recreate a target image of dimensions $H_V \times W_V$, we need to evaluate N points $\{p_{u,v}^z\}_{z=1}^N$ along each ray $r_{u,v}$ with direction $d_{u,v}$, where $u \in \llbracket 1, H_V \rrbracket$ and $v \in \llbracket 1, W_V \rrbracket$:

$$d_{u,v} = R_t K^{-1} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} + T_t, \quad p_{u,v}^z = T_t + t_z \frac{d_{u,v}}{\|d_{u,v}\|}, \quad (2)$$

where $t_z \sim \mathcal{U}[z_n + \frac{z-1}{N}(z_f - z_n), z_n + \frac{z}{N}(z_f - z_n)]$ following [3], z_n and z_f being the depth near and far bounds of the visual frustum. K is the intrinsic camera matrix. The target volume $F_{t/i}$ is obtained then as the resampling of input volume F_i with trilinear interpolation, using points $\{p_{u,v}^z\}$ aligned rigidly to the input camera coordinate frame:

$$F_{t/i}(u, v, z) = F_i(R_i^\top(p_{u,v}^z - T_i)), \quad (3)$$

where $F_{t/i} \in \mathbb{R}^{H_V \times W_V \times N \times C}$ and $\{R_i, T_i\}$ is the input camera pose. In practice, we normalize the aligned points' coordinates prior to sampling as F_i is assumed to represent features in the input view normalized device coordinate (NDC) space. We use the NDC parametrization for an optimal spatial exploitation of the input feature volume F_i and generalization across objects and scenes with different scales and datasets with different camera settings (e.g. intrinsics, z_n , z_f).

3.3 Feature Aggregation

As different input views provide different information about the observed scene, we merge subsequently the 3D features obtained from the various inputs. We note that all target feature volumes $\{F_{t/i}^k\}_k$ provided by input images $\{I_i^k\}_k$ are represented in the same target view camera coordinate frame. A naive merging strategy would be to simply average these volumes element-wise. However, for a given 3D location in the target view frustum, different input views contribute appearance information with varying confidence, based on the visibility/occlusion of this spatial location in the input views. In order to emulate this principle, and inspired by attention mechanisms, we propose to learn a 3D confidence measure per input view in the form of a weight volume $W_i \in \mathbb{R}^{H_V \times W_V \times D}$. This volume is obtained as one of the channels of the input volume features $W_i = F_i(1)$ (i.e. $W_{t/i} = F_{t/i}(1)$). As this confidence volume depends naturally on the input image and the relative camera pose of the target with respect to the input, similarly to [31], we append these relative poses to the input image pixel values as additional input to the encoder E . After resampling the input features $\{F_i^k\}_k$ into the target ones $\{F_{t/i}^k\}_k$, we use the resampled weights $\{W_{t/i}^k\}_k$ normalized with Softmax across the input views to compute a weighted average of the target volumes:

$$F_t = \sum_k \text{Softmax}_k(W_{t/i}^k) F_{t/i}^k(\llbracket 1, C \rrbracket), \quad (4)$$

where index k is over the number of input views, and $F_t \in \mathbb{R}^{H_V \times W_V \times N \times C-1}$. Let us recall that this tensor represents features of N points, within $[z_n, z_f]$ depth-wise, for all rays associated with a target image of dimension $H_V \times W_V$. This aggregation allows our method to use an arbitrary number of input views at both training and testing. In the single input case, we note that $F_t = F_{t/i}(\llbracket 1, C \rrbracket)$.

3.4 Feature Rendering

Following volumetric rendering [48,3], we generate a target feature image $\tilde{\mathcal{F}}$ for a given target view differentially using α -compositing of the target feature volume F_t along the depth dimension. To this end, we assume one of the target feature channels to represent volume density [48] $\sigma = F_t(1) \in \mathbb{R}^{H_V \times W_V \times D}$. We recall that the dimensions of tensor F_t span the pixels of the target feature resolution $H_v \times W_v$ in the first two dimensions, and N points sampled along each ray for the third dimension. The rendered target feature image then writes:

$$\tilde{\mathcal{F}} = \sum_{z=1}^N T_z \alpha_z F_t(\llbracket 1, C-1 \rrbracket), \quad (5)$$

$$T_z = e^{-\sum_{j=1}^{z-1} \sigma(j)\delta_j}, \quad \alpha_z = 1 - e^{\sigma(z)\delta_z}, \quad (6)$$

where T represents transmittance, $\delta_z = t_{z+1} - t_z$ and $\tilde{\mathcal{F}} \in \mathbb{R}^{H_V \times W_V \times C-2}$. In order to reduce the memory cost and increase the rendering speed of our method, the size of the rendered feature image is chosen to be lower than the size of the target image resolution, i.e. $H_V = H/4$ and $W_V = W/4$.

3.5 Neural Light Field

At this stage, the convolutional rendered features produce a low resolution feature image representative of all rays making up the target view. We propose to learn a light field function f , which performs both upsampling and refinement of the result of the convolutional first stage of our method. This implicit neural network maps rays of the target image to their colors, while being conditioned on ray features extracted from the convolutional rendered features.

Given a ray $r_{u,v}$ with direction $d_{u,v}$ corresponding to the target image pixel coordinates (u, v) , with $(u, v) \in \llbracket 1, H \rrbracket \times \llbracket 1, W \rrbracket$, we encode rays using Plücker coordinates similarly to Sitzmann et al. [11]:

$$r_{u,v} = \frac{(d_{u,v}, T_t \times d_{u,v})}{\|d_{u,v}\|}, \quad (7)$$

where $r_{u,v} \in \mathbb{R}^6$. This representation ensures a unique ray encoding when the origin T_t moves along direction $d_{u,v}$. We recall that the expression of $d_{u,v}$ as a function of the target camera pose $\{R_t, T_t\}$ can be found in equation 2.

The feature $\mathcal{F}_{u,v}$ of a ray $r_{u,v}$ at the final image resolution $H \times W$ is obtained from the lower resolution rendered feature image $\tilde{\mathcal{F}} \in \mathbb{R}^{H_v \times W_v \times C-2}$ through a learned upsampling. Specifically, the rendered feature image undergoes two successive 2D convolutions and up samplings to produce a feature image at the desired resolution $\mathcal{F} \in \mathbb{R}^{W \times H \times C-2}$. The final target RGB image $I_t = \{c_{u,v}\}_{u \in \llbracket 1, H \rrbracket, v \in \llbracket 1, W \rrbracket}$ is predicted then from the concatenation of the ray coordinate and its feature with an MLP accordingly:

$$c_{u,v} = f(r_{u,v}, \mathcal{F}_{u,v}). \quad (8)$$

We refer the reader to the supplementary material for more details about the architecture of network f . Notice that while convolution equipped NeRF [3] methods (e.g. pixelNeRF [7] MVSNeRF [10] GRF[45]) require querying $H \times W \times N$ 3D points through their implicit neural radiance fields, our light field network only needs to evaluate $H \times W$ rays, which enables our method to train potentially faster, and render orders of magnitude faster compared to e.g. pixelNeRF (see Table 2).

3.6 Training Objective

Our model is fully differentiable and trained end-to-end. We optimize the parameters of the model, namely the convolutional network E and the light field network f jointly, by back-propagating a combination of a fine loss L_r and two coarse losses \tilde{L}_r and \tilde{L}_d :

$$L = L_r + \tilde{L}_r + \tilde{L}_d. \quad (9)$$

L_r is a L2 reconstruction loss between the final image I_t predicted by the light field network and the ground-truth I_t^{gt} :

$$L_r = \|I_t - I_t^{gt}\|_2^2. \quad (10)$$

We regularize the first convolutional stage of our method through a coarse L2 reconstruction loss \tilde{L}_r . After rendering the low resolution image feature $\tilde{\mathcal{F}}$, we constrain its first 3 channels to produce the RGB colors of the target ground truth image:

$$\tilde{L}_r = \|\tilde{I}_t - \tilde{I}_t^{gt}\|_2^2, \quad (11)$$

where $\tilde{I}_t = \tilde{\mathcal{F}}([1, 3])$, and \tilde{I}_t^{gt} is the ground truth image downsampled to resolution $H_V \times W_V$.

Additionally, we regularize the gradient of the low resolution depth image \tilde{d}_t rendered from the density volume σ of the first stage. Similarly to [49], we weight this cost with an edge-aware term using the ground-truth image gradient:

$$\tilde{L}_d = \frac{1}{H_V \times W_V} \sum_{u,v} |\partial_u \tilde{d}_t| e^{-\|\partial_u \tilde{I}_t^{gt}\|} + |\partial_v \tilde{d}_t| e^{-\|\partial_v \tilde{I}_t^{gt}\|}. \quad (12)$$

The depth image \tilde{d}_t can be expressed as a function of transmittance T and α values as follows:

$$\tilde{d}_t = \frac{1}{\sum_{z=1}^N T_z \alpha_z} \sum_{z=1}^N T_z \alpha_z t_z. \quad (13)$$

We note that the expressions of T and α are detailed in equation 6.

4 Experiments

We demonstrate the capability of our method to generate novel views from sparse input views using standard benchmarks. Specifically, we evaluate the ability of our network to reconstruct 360-degree novel views of objects unseen at training in the synthetic ShapeNet dataset [4]. We additionally evaluate our ability to infer novel views of real world scenes using the more challenging DTU multi-view stereo dataset [13].

4.1 Implementation Details

We implemented our method with the PyTorch [50] framework on a Quadro RTX 5000 gpu. We optimize with the Adam [51] solver using learning rate 10^{-4} in training and 10^{-5} in fine-tuning. The depth of the convolutional feature volume is set to $D = 32$, and the number of channels $C = 32$. For the MLP of the light field network, we use 5 layers with a hidden dimension of 256, similarly to NeRF [3]. For volumetric feature resampling, we use the coarse and fine sampling strategy similarly to previous work (e.g. [3,7]), with $N = 64$ coarse samples and $N = 32$ fine samples. We show detailed network structure in the supplementary material.

4.2 Generalization on Synthetic Data

Dataset Following the settings in GRF [45] and PixelNeRF [7], we evaluate first our synthesis results on the car and chair classes of dataset ShapeNet-V2 [4]. Precisely, the cars amount to 2151 training objects, 352 validation objects and 704 testing objects, while the chairs count 4612 training objects, 662 validation objects and 1317 testing objects. Hence the testing objects are not seen in training. In the training split, each object has 50 images with size 128×128 . For testing, there are 251 views per object. Our model takes 1 or 2 fixed views as input and infers novel views for evaluation.

Method	PSNR(Cars) \uparrow		SSIM(Cars) \uparrow		PSNR(Chairs) \uparrow		SSIM(Chairs) \uparrow	
	1-view	2-view	1-view	2-view	1-view	2-view	1-view	2-view
SRN[4]	20.72	22.94	0.85	0.88	22.89	24.48	0.91	0.92
LFN[11]	22.42	–	0.89	–	22.26	–	0.90	–
TCO[27]	18.15	18.41	0.79	0.80	21.27	21.33	0.88	0.88
WRL[28]	16.89	17.20	0.77	0.78	22.11	22.28	0.90	0.90
dGQN[52]	18.19	18.79	0.78	0.79	21.59	22.36	0.87	0.89
PixelNeRF[7]	23.17	25.66	0.90	0.94	23.72	26.20	0.91	0.94
GRF[45]	20.33	22.34	0.82	0.86	21.25	22.65	0.86	0.88
ENR[35]	22.26	–	–	–	22.83	–	–	–
Ours	22.31	23.82	0.87	0.91	22.52	24.10	0.90	0.92

Table 1: Comparison of the average PSNR and SSIM of reconstructed images in the ShapeNet-V2 [4] dataset. The higher the better for both PSNR and SSIM. Red, orange and blue represent the first, second and third ranking methods respectively. SRN[4] and LFN[11] require test time optimization.

	pixelNeRF[7]	ours
clock time for a 128×128 image in seconds	6.23	0.06

Table 2: Comparison of rendering complexity. All clock times are collected on a Quadro RTX 5000 gpu. Our rendering is 100 times faster than PixelNeRF [7].

Table 1 shows a quantitative comparison of our method with the recent state-of-the-art in few-shot view synthesis. We report the peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) reconstruction metrics. We relay the numbers for methods TCO[27], WRL[28], dGQN[52] SRN[4] and GRF [45] as they were reported in [45]. We report the numbers of ENR[35] and PixelNeRF [7] from [7], and the numbers for LFN[11] from their paper. Figure 4 shows a qualitative comparison to a few of these methods. We obtain the visualizations for PixelNeRF [7] and LFN[11] using their publicly available codes and models. The teaser figure 1 shows additional 360-degree novel view synthesis results, and we provide further visual results in the supplementary material.

The results confirm that 2D based image to image novel view methods (e.g. TCO) are outperformed by the 3D aware ones (e.g. ENR, SRN, etc). Furthermore, 3D aware methods that use implicit 3D representations (e.g. PixelNeRF,

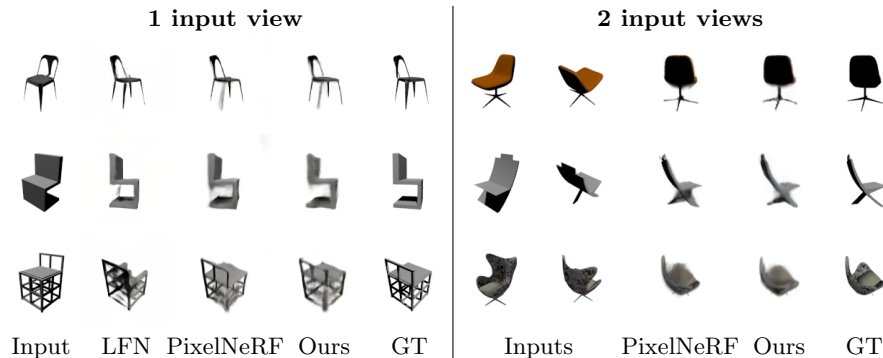


Fig. 3: Qualitative comparison of novel view synthesis of unseen chairs from a single and 2 input views on ShapeNet-V2 [4] (More results in **supplementary material**). LFN [11] requires test times optimization.

GRF, SRN) outperform generally their counterparts relying on explicit 3D latents (e.g. ENR). Our method is hybrid, in that it uses an explicit 3D latent, combined with a 2D implicit representation.

Numerically, our method is overall second to PixelNeRF, while being two orders of magnitude faster at rendering. As we show in table 2, our method requires only 0.06 seconds to infer a full 128×128 image, while PixelNeRF takes 6.23 seconds using the implementation provided by the authors. As shown in figure 4, our method provides comparable reconstruction quality as well, while encountering similar difficulties for challenging views.

While our performance is generally close to LFN across the benchmark, we note that LFN requires auto-decoding test time optimization. It also requires training hypernetworks, which are prohibitively expensive in compute and memory, thereby limiting the resolution of the reconstructed images. This hinders in turn the applicability of this method to real datasets with images larger than 128×128 . Conversely, we demonstrate the ability of our method to model complex real scenes with bigger images using moderate computational resources, while providing optimization-free single forward pass prediction. We note that SRN requires test-time optimization as well.

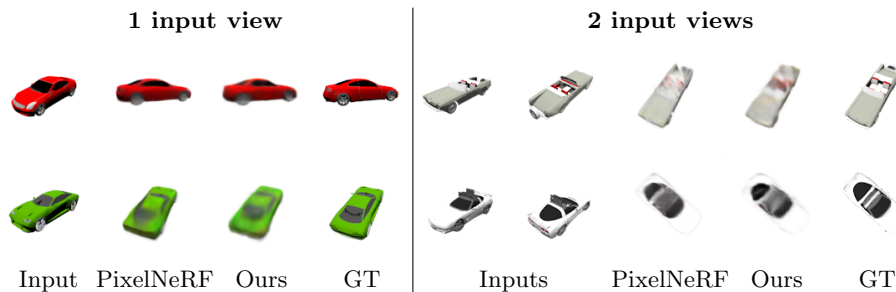


Fig. 4: Qualitative comparison of novel view synthesis of unseen cars from a single and 2 input views on ShapeNet-V2 [4] (More results in **supp. material**).

4.3 Generalization on Real Data

Dataset We demonstrate that our method is capable of reconstructing novel views for real world scenes unseen at training using the DTU dataset [13]. Following the PixelNeRF [7] experimental settings, the data is split into 88 training scenes and 16 testing scenes, each scene including 49 images with resolution 300×400 . We note that this is a challenging scenario. In fact, the training scenes are limited, and the training and testing scenes do not share any semantic similarities as can be seen in figure 5. The lighting and backgrounds are also inconsistent between the scenes. Hence, this is a few-shot novel view synthesis task that demands considerable scene category generalization as well.

Method	PSNR \uparrow		SSIM \uparrow		LPIPS \downarrow	
	3-view	6-view	3-view	6-view	3-view	6-view
SRF[53]	15.84	17.77	0.532	0.616	0.482	0.401
PixelNeRF[7]	19.33	20.43	0.695	0.732	0.387	0.361
MVSNeRF[10]	16.33	18.26	0.602	0.695	0.385	0.321
Ours	17.31	19.20	0.611	0.655	0.433	0.398

Table 3: Comparison of the average PSNR, SSIM and LPIPS of reconstructed images in the DTU [13] dataset **without test time optimization**. The higher the better for both PSNR and SSIM. The lower the better for LPIPS. Red, orange and blue represent the first, second and third ranking methods respectively.

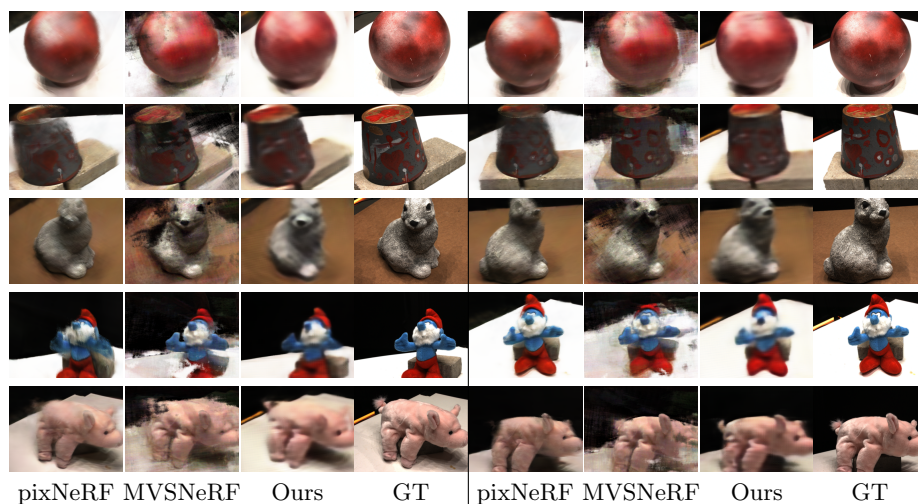


Fig. 5: Qualitative comparison of novel view synthesis of unseen scenes without test time optimization from 6 input views on the DTU dataset [13].

Table 3 shows a quantitative comparison of our method with the recent state-of-the-art in optimization-free few-shot view synthesis. We report the peak signal-to-noise ratio (PSNR), structural similarity (SSIM) and learned perceptual image patch similarity (LPIPS) reconstruction metrics, for the same 3 and

6 view inputs averaged across the same testing scenes. For a fair comparison, we report the performance of PixelNeRF [7] from their paper, and the numbers of methods MVSNeRF [10] and SRF [53] from RegNeRF [39], as the authors in the latter reproduce the performance of these methods in pixelNeRF’s DTU setup. Figure 5 shows a qualitative comparison between our method and methods pixelNeRF [7] and MVSNeRF [10] on synthesized views from testing scenes given the same inputs. We produce the results of PixelNeRF using their publicly available code and DTU model. For MVSNeRF, as their original model was trained on a different DTU setup, we finetune their model on the pixelNeRF DTU setup similarly to RegNeRF [39].

Although PixelNeRF and MVSNeRF are based on implicit radiance fields and hence require more evaluations and time for rendering, our implicit light field based method provides competitive performances in comparison. In the single forward prediction setting, our method is overall second to PixelNeRF in PSNR, while providing a 100 times faster inference (see table 2). As illustrated in the visual comparison in figure 5, we manage to reproduce the shape and appearance of the scene to a good extent and also recover from some of the competitions’ failures. Our method appears to be better in fact at preserving the coarser structure of the scene. Specifically, some elements of the ground truth that we manage to reproduce are not recovered by the competition, such as the pig tail in the last row, the background table’s yellow lines and the rabbit’s eyes. Although we found MVSNeRF encounters multiple failures compared to pixelNeRF, it is apparent that NeRF base methods (pixelNeRF and MVSNeRF) are able to produce some relatively higher frequency details, albeit at a considerable higher rendering cost.

Method	PSNR \uparrow		SSIM \uparrow		LPIPS \downarrow	
	3-view	6-view	3-view	6-view	3-view	6-view
mip-NeRF[40]	7.64	14.33	0.227	0.568	0.655	0.394
DietNeRF [38]	10.01	18.70	0.354	0.668	0.574	0.336
RegNeRF [39]	15.33	19.10	0.621	0.757	0.341	0.233
SRF [53]	16.06	18.69	0.550	0.657	0.431	0.353
PixelNeRF[7]	17.38	21.52	0.548	0.670	0.456	0.351
MVSNeRF[10]	16.26	18.22	0.601	0.694	0.384	0.319
Ours	17.72	19.56	0.626	0.671	0.412	0.375

Table 4: Comparison of the average PSNR, SSIM and LPIPS of reconstructed images in the DTU [13] dataset **with test time optimization**. The higher the better for both PSNR and SSIM. The lower the better for LPIPS. Red, orange and blue represent the first, second and third ranking methods respectively.

Table 4 shows a quantitative comparison of our method with the recent few-shot novel view synthesis state-of-the-art with test time optimization. We report PSNR, SSIM and LPIPS for the same 3 and 6 input views averaged over the same testing objects. All the method’s numbers on the PixelNeRF DTU setup are reported as reproduced in RegNeRF [39]. Methods mip-NeRF [40], DietNeRF [38] and RegNeRF [39] are optimized per scene only, while PixelNeRF [7],

MVSNeRF [10] and ours are trained on the DTU training set then finetuned per scene. Following the experimental setting in RegNeRF, only the input views were used for fine-tuning. Similarly to the finetuning of MVSNeRF and PixelNeRF, we reduce the learning rate from 10^{-4} to 10^{-5} and constrain the finetuning within 10k iterations for better performance. Figure 6 shows a qualitative comparison to MVSNeRF and PixelNeRF given 6 input views after finetuning.

In the 3 input view case, we outperform all methods in the PSNR and SSIM metrics. We obtain overall comparable performances with generalizable (PixelNeRF and MVSNeRF) and single scene optimization (RegNeRF) NeRFs. Figure 6 shows that we can achieve relatively comparable results to the encoder endowed NeRF approaches after optimization. We recall again that competition methods here require renderings that are orders of magnitude slower than ours.

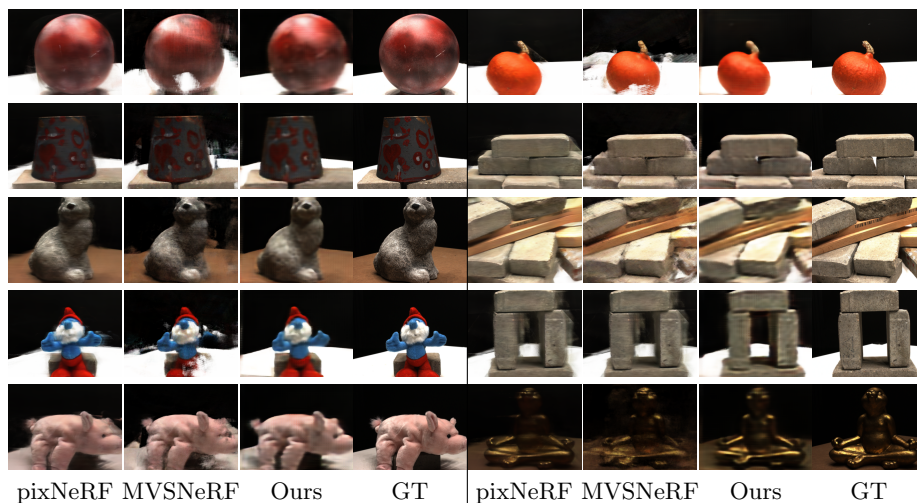


Fig. 6: Qualitative comparison of novel view synthesis with test time optimization using 6 input views on the DTU dataset [13].

4.4 Ablation

We propose here an ablative analysis for our method on the DTU [13] and shapeNet-V2 [4] datasets. Specifically, we disable the light field function (ours w/o light field), and we render the final image directly from the target view aligned convolutional feature volume. We also reproduce our method without using the ray coordinates (ours w/o ray coordinates).

Table 5 shows numerical comparisons for 3 and 6 input views on DTU, and figure 7 shows qualitative comparisons for 6 input views on DTU, and 1 input view on ShapeNet-V2. Our final method improves on its 3D aware convolutional baseline both numerically and visually. In particular, we note how the light field network with ray encoding reduces ghosting artifacts and provides finer details.

Method	PSNR \uparrow		SSIM \uparrow		LPIPS \downarrow	
	3-view	6-view	3-view	6-view	3-view	6-view
Ours w/o light field	16.61	17.30	0.575	0.603	0.462	0.402
Ours w/o ray coordinates	16.77	17.42	0.604	0.636	0.455	0.403
Ours	17.31	19.20	0.611	0.655	0.433	0.398

Table 5: Quantitative ablation of our method on the DTU dataset [13]. The higher the better for both PSNR and SSIM. The lower the better for LPIPS.



Fig. 7: Qualitative ablation of our method on unseen DTU [13] scenes (6 input views) and unseen ShapeNet-V2[4] cars (1 input view).

5 Limitations

Whilst our method offers an efficient rendering (100 times faster than PixelNeRF), it still has difficulties in reproducing the highest level of details in real large images as the NeRF based methods. We believe this is due to the reduced resolution of our feature volume and our coarse feature rendering, which conversely contribute towards reducing the memory footprint. As future work, we will attempt to overcome this limitation while maintaining memory efficiency. Furthermore, the explicit feature volume limits the applicability of our method on datasets with larger images such as LLFF [16]. We will thus consider intrinsic 3D feature representations next.

6 Conclusion

We proposed a method for generating novel views from few input calibrated images with a single forward pass prediction deep neural network. We learn an implicit neural light field function that models ray colors directly. In comparison to [11], we proposed a more efficient local ray conditioning, and an optimization free inference. Our method combines the advantages of 3D aware convolutional approaches and implicit representations, and requires only image data in training. We demonstrated our method successfully on synthetic and real benchmarks for few-shot novel view synthesis. Our method outperforms the convolutional baselines (see table 1) and provides competitive performances compared to locally conditioned radiance fields (e.g. PixelNeRF [7]), while being a 100 times faster to render.

References

1. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2019) 165–174 [2](#)
2. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2019) 4460–4470 [2](#)
3. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: European conference on computer vision, Springer (2020) 405–421 [2, 3, 5, 6, 7, 8, 19](#)
4. Sitzmann, V., Zollhöfer, M., Wetzstein, G.: Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems* **32** (2019) [2, 3, 8, 9, 10, 13, 14, 18, 20, 21, 22](#)
5. Peng, S., Niemeyer, M., Mescheder, L., Pollefeys, M., Geiger, A.: Convolutional occupancy networks. In: ECCV. (2020) [2](#)
6. Chibane, J., Alldieck, T., Pons-Moll, G.: Implicit functions in feature space for 3d shape reconstruction and completion. In: CVPR. (2020) [2](#)
7. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2021) 4578–4587 [2, 3, 7, 8, 9, 11, 12, 14, 18](#)
8. Sitzmann, V., Chan, E.R., Tucker, R., Snavely, N., Wetzstein, G.: Metasdf: Meta-learning signed distance functions. In: NeurIPS. (2020) [2](#)
9. Ouasfi, A., Boukhayma, A.: Few’zero level set’-shot learning of shape signed distance functions in feature space. arXiv preprint arXiv:2207.04161 (2022) [2](#)
10. Chen, A., Xu, Z., Zhao, F., Zhang, X., Xiang, F., Yu, J., Su, H.: Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 14124–14133 [2, 3, 4, 7, 11, 12, 13](#)
11. Sitzmann, V., Rezkikov, S., Freeman, B., Tenenbaum, J., Durand, F.: Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems* **34** (2021) [2, 3, 7, 9, 10, 14, 18, 22](#)
12. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015) [2](#)
13. Jensen, R., Dahl, A., Vogiatzis, G., Tola, E., Aanæs, H.: Large scale multi-view stereopsis evaluation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2014) 406–413 [2, 3, 8, 11, 12, 13, 14, 18](#)
14. Tucker, R., Snavely, N.: Single-view view synthesis with multiplane images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2020) 551–560 [2](#)
15. Shih, M.L., Su, S.Y., Kopf, J., Huang, J.B.: 3d photography using context-aware layered depth inpainting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2020) 8028–8038 [2](#)
16. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* **38** (2019) 1–14 [2, 14](#)

17. Niklaus, S., Mai, L., Yang, J., Liu, F.: 3d ken burns effect from a single image. *ACM Transactions on Graphics (ToG)* **38** (2019) 1–15 [2](#)
18. Choi, I., Gallo, O., Troccoli, A., Kim, M.H., Kautz, J.: Extreme view synthesis. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. (2019) 7781–7790 [2](#)
19. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In: *NeurIPS*. (2021) [3](#)
20. Kellnhofer, P., Jebe, L.C., Jones, A., Spicer, R., Pulli, K., Wetzstein, G.: Neural lumigraph rendering. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. (2021) 4287–4297 [3](#)
21. Rematas, K., Martin-Brualla, R., Ferrari, V.: Sharf: Shape-conditioned radiance fields from a single view. *arXiv preprint arXiv:2102.08860* (2021) [3](#)
22. Wei, Y., Liu, S., Rao, Y., Zhao, W., Lu, J., Zhou, J.: Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. (2021) 5610–5619 [3](#)
23. Riegler, G., Koltun, V.: Free view synthesis. In: *European Conference on Computer Vision*, Springer (2020) 623–640 [3](#)
24. Riegler, G., Koltun, V.: Stable view synthesis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. (2021) 12216–12225 [3](#)
25. Thies, J., Zollhöfer, M., Nießner, M.: Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)* **38** (2019) 1–12 [3](#)
26. Aliev, K.A., Sevastopolsky, A., Kolos, M., Ulyanov, D., Lempitsky, V.: Neural point-based graphics. In: *European Conference on Computer Vision*, Springer (2020) 696–712 [3](#)
27. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Single-view to multi-view: Reconstructing unseen views with a convolutional network. *CoRR* abs/1511.06702 **1** (2015) [2](#) [3](#), [9](#)
28. Worrall, D.E., Garbin, S.J., Turmukhambetov, D., Brostow, G.J.: Interpretable transformations with encoder-decoder networks. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2017) 5726–5735 [3](#), [9](#)
29. Zhou, T., Tulsiani, S., Sun, W., Malik, J., Efros, A.A.: View synthesis by appearance flow. In: *European conference on computer vision*, Springer (2016) 286–301 [3](#)
30. Park, E., Yang, J., Yumer, E., Ceylan, D., Berg, A.C.: Transformation-grounded image generation network for novel 3d view synthesis. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. (2017) 3500–3509 [3](#)
31. Sun, S.H., Huh, M., Liao, Y.H., Zhang, N., Lim, J.J.: Multi-view to novel view: Synthesizing novel views with self-learned confidence. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. (2018) 155–171 [3](#), [6](#)
32. Lombardi, S., Simon, T., Saragih, J., Schwartz, G., Lehrmann, A., Sheikh, Y.: Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751* (2019) [3](#)
33. Wiles, O., Gkioxari, G., Szeliski, R., Johnson, J.: Synsin: End-to-end view synthesis from a single image. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. (2020) 7467–7477 [3](#)
34. Olszewski, K., Tulyakov, S., Woodford, O., Li, H., Luo, L.: Transformable bottleneck networks. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. (2019) 7648–7657 [3](#), [4](#)

35. Dupont, E., Martin, M.B., Colburn, A., Sankar, A., Susskind, J., Shan, Q.: Equivariant neural rendering. In: International Conference on Machine Learning, PMLR (2020) 2761–2770 [3](#), [4](#), [9](#)
36. Hu, R., Ravi, N., Berg, A.C., Pathak, D.: Worldsheet: Wrapping the world in a 3d sheet for view synthesis from a single image. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 12528–12537 [3](#)
37. Chen, X., Song, J., Hilliges, O.: Monocular neural image based rendering with continuous view control. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2019) 4090–4100 [3](#)
38. Jain, A., Tancik, M., Abbeel, P.: Putting nerf on a diet: Semantically consistent few-shot view synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 5885–5894 [3](#), [12](#)
39. Niemeyer, M., Barron, J.T., Mildenhall, B., Sajjadi, M.S., Geiger, A., Radwan, N.: Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In: CVPR. (2022) [3](#), [12](#), [18](#)
40. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 5855–5864 [3](#), [12](#)
41. Liu, L., Gu, J., Zaw Lin, K., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. *Advances in Neural Information Processing Systems* **33** (2020) 15651–15663 [3](#)
42. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 5752–5761 [3](#)
43. Garbin, S.J., Kowalski, M., Johnson, M., Shotton, J., Valentin, J.: Fastnerf: High-fidelity neural rendering at 200fps. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 14346–14355 [3](#)
44. Reiser, C., Peng, S., Liao, Y., Geiger, A.: Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 14335–14345 [3](#)
45. Trevithick, A., Yang, B.: Grf: Learning a general radiance field for 3d representation and rendering. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 15182–15192 [3](#), [7](#), [9](#)
46. Wang, Q., Wang, Z., Genova, K., Srinivasan, P.P., Zhou, H., Barron, J.T., Martin-Brualla, R., Snavely, N., Funkhouser, T.: Ibrnet: Learning multi-view image-based rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2021) 4690–4699 [3](#)
47. Hedman, P., Srinivasan, P.P., Mildenhall, B., Barron, J.T., Debevec, P.: Baking neural radiance fields for real-time view synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. (2021) 5875–5884 [3](#)
48. Kajiya, J.T., Von Herzen, B.P.: Ray tracing volume densities. In: Proceedings of the 11th Annual Conference on Computer Graphics and Interactive Techniques. SIGGRAPH '84 (1984) 165–174 [5](#), [6](#)
49. Heise, P., Klose, S., Jensen, B., Knoll, A.: Pm-huber: Patchmatch with huber regularization for stereo matching. In: Proceedings of the IEEE International Conference on Computer Vision. (2013) 2360–2367 [8](#)
50. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019) [8](#)

51. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) 8
52. Eslami, S.A., Jimenez Rezende, D., Besse, F., Viola, F., Morcos, A.S., Garnelo, M., Ruderman, A., Rusu, A.A., Danihelka, I., Gregor, K., et al.: Neural scene representation and rendering. *Science* **360** (2018) 1204–1210 9
53. Chibane, J., Bansal, A., Lazova, V., Pons-Moll, G.: Stereo radiance fields (srf): Learning view synthesis for sparse views of novel scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2021) 7911–7920 11, 12

Supplementary Material

Network architecture

Table 6 describes the detailed architecture of our convolutional network E introduced in Section 3.1 of the main submission. Table 7 shows the detailed architecture of our network f introduced in Section 3.5 in the main submission.

DTU [13] evaluation protocol

We follow the evaluation setup in PixelNeRF [7,39] and use the following scans as the test set: 8, 21, 30, 31, 34, 38, 40, 41, 45, 55, 63, 82, 103, 110, 114. The following images are used as input: 25, 22, 28, 40, 44, 48, 0, 8, 13. For the 3 and 6 input scenarios, we use the first 3/6 images from the list. All remaining images are used for evaluation except wrong exposure images 3, 4, 5, 6, 7, 16, 17, 18, 19, 20, 21, 36, 37, 38, 39. The input image resolution is 300×400 .

Additional qualitative results

Results on ShapeNet-V2 [4] Figures 8, 9, 10, 11 show additional novel view synthesis qualitative results of our method using a single and 2 input views.

Comparison to LFN [11] Figure 12 shows additional comparisons to Sitzmann et al. [11] on the chair class of ShapeNet-V2 [4].

Input Shape	Output shape	Operation	
$(3, H, W)$	$(52, H, W)$	1×1 Conv	Image
$(12, H, W)$	$(12, H, W)$	1×1 Conv	Relative Pose
$(52 + 12, H, W)$	$(64, H, W)$	1×1 Conv	
$(64, H, W)$	$(64, H, W)$	$2 \times$ ResBlock	
$(64, H, W)$	$(128, H/2, W/2)$	4×4 Conv,Stride2	
$(128, H/2, W/2)$	$(128, H/2, W/2)$	$1 \times$ ResBlock	
$(128, H/2, W/2)$	$(128, H/4, W/4)$	4×4 Conv,Stride2	2D Conv
$(128, H/4, W/4)$	$(128, H/4, W/4)$	$1 \times$ ResBlock	
$(128, H/4, W/4)$	$(256, H/8, W/8)$	4×4 Conv,Stride2	
$(256, H/8, W/8)$	$(256, H/8, W/8)$	$1 \times$ ResBlock	
$(256, H/8, W/8)$	$(128, H/4, W/4)$	4×4 Conv.T,Stride2	
$(128, H/4, W/4)$	$(128, H/4, W/4)$	$2 \times$ ResBlock	
$(128, H/4, W/4)$	$(256, H/4, W/4)$	1×1 Conv	
$(256, H/4, W/4)$	$(512, H/4, W/4)$	1×1 Conv	2D to 3D
$(512, H/4, W/4)$	$(2048, H/4, W/4)$	1×1 Conv	
$(2048, H/4, W/4)$	$(32, 64, H/4, W/4)$	Reshape	
$(32, 64, H/4, W/4)$	$(32, 64, H/4, W/4)$	1×1 Conv	
$(32, 64, H/4, W/4)$	$(32, 64, H/4, W/4)$	$2 \times$ ResBlock	
$(32, 64, H/4, W/4)$	$(64, 32, H/8, W/8)$	4×4 Conv,Stride2	3D Conv
$(64, 32, H/8, W/8)$	$(64, 32, H/8, W/8)$	$2 \times$ ResBlock	
$(64, 32, H/8, W/8)$	$(32, 64, H/4, W/4)$	4×4 Conv.T,Stride2	
$(32, 64, H/4, W/4)$	$(32, 64, H/4, W/4)$	$2 \times$ ResBlock	
$(31, 64, H/4, W/4)$	$(31, H/4, W/4)$	-	Rendering
$(30, H/4, W/4)$	$(30, H/2, W/2)$	1×1 Conv,Upsampling	Upsampling
$(30, H/2, W/2)$	$(30, H, W)$	1×1 Conv,Upsampling	

Table 6: The architecture of network E in Section 3.1 of the main submission.

Layer	Channels	Inputs
LR_0	30/256	$\mathcal{F}_{u,v}$
PE	6/78	$r_{u,v}$
LR_1	78/256	PE
LR_2	256/256	$LR_1 \odot LR_0$
LR_3	256/256	$LR_2 \odot LR_0$
LR_4	256/256	$LR_3 \odot LR_0$
LR_5	256/256	$LR_4 \odot LR_0$
LR_6	256/256	$LR_5 \odot LR_0$
$c_{u,v}$	256/3	LR_6

Table 7: The MLP in Section 3.5 of the main submission. $r_{u,v}$ is the Plücker coordinate of the ray. $\mathcal{F}_{u,v}$ is the feature of the ray. $c_{u,v}$ is the target pixel color. PE is the positional encoding in [3]. \odot is the element-wise product. All layers are linear with Relu activation, except the last layer which uses a Sigmoid.

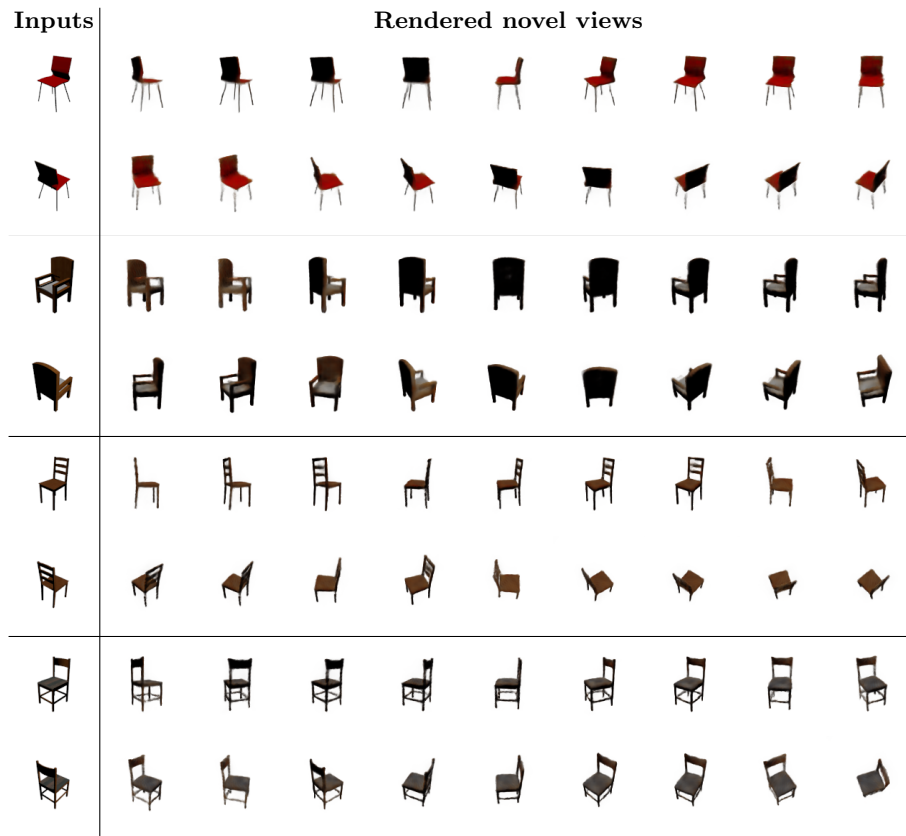


Fig. 8: Novel view synthesis of chairs from ShapeNet-V2 [4] given 2 input views using our method.

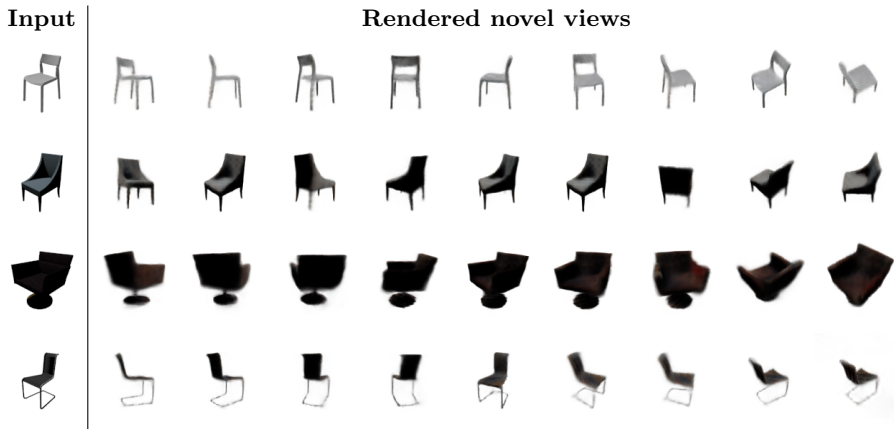


Fig. 9: Novel view synthesis of chairs from ShapeNet-V2 [4] given a single input view using our method.

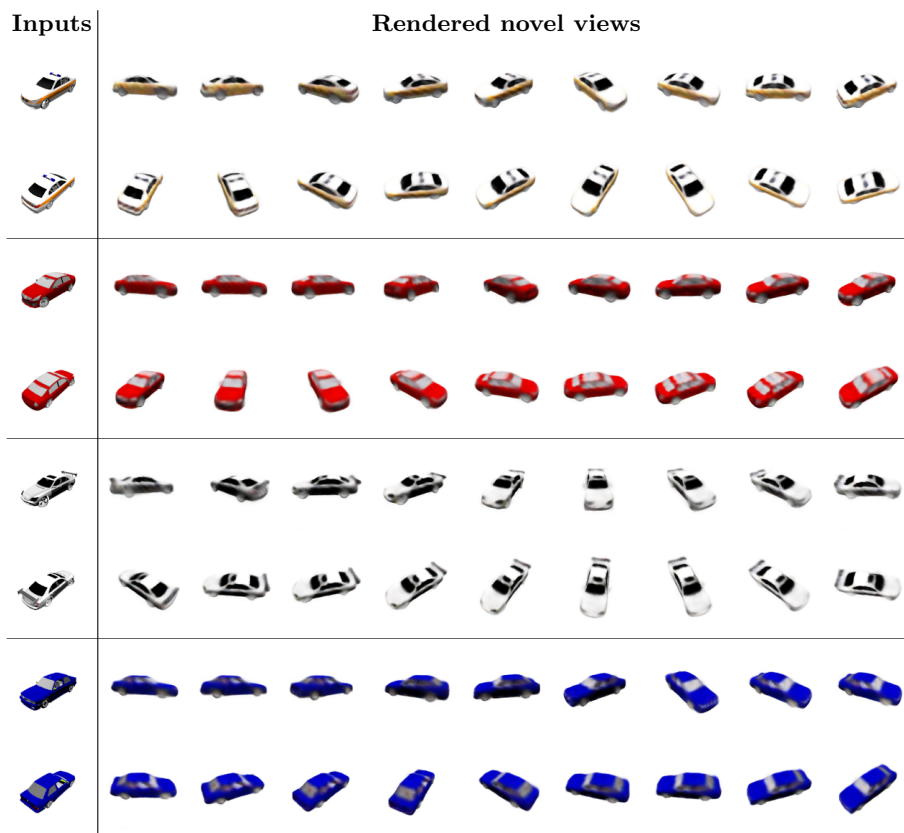


Fig. 10: Novel view synthesis of cars from ShapeNet-V2 [4] given 2 input views using our method.

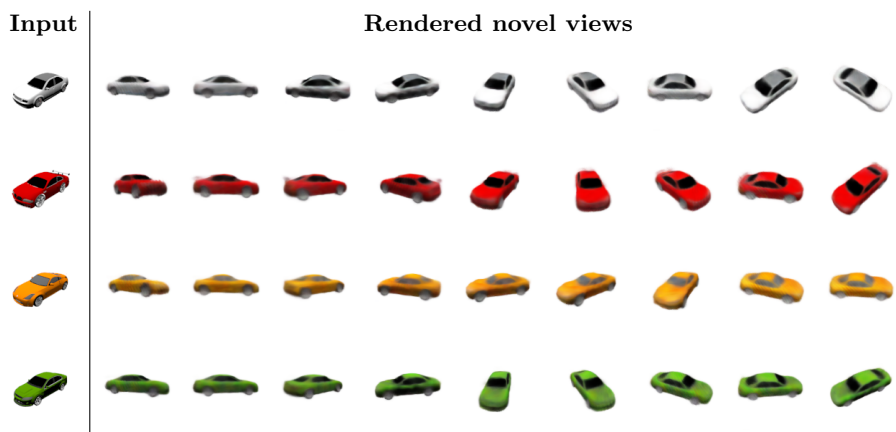


Fig. 11: Novel view synthesis of cars from ShapeNet-V2 [4] given a single input view using our method.

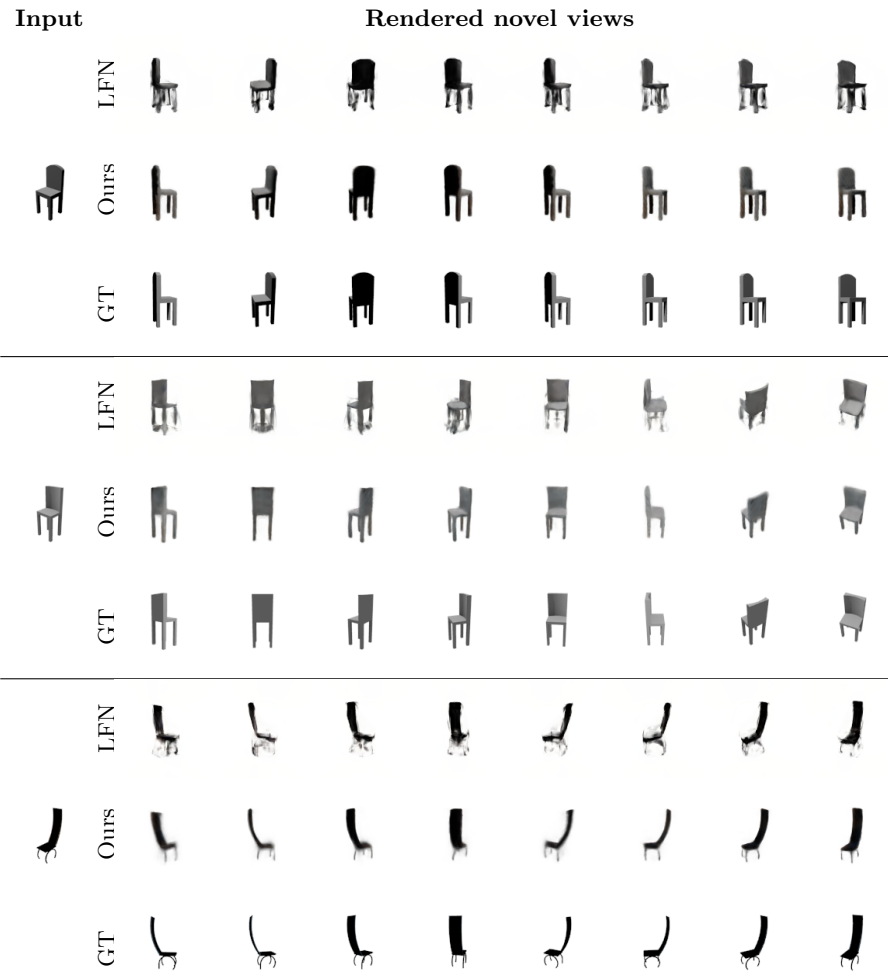


Fig. 12: Qualitative comparison to [11] on novel view synthesis of chairs from ShapeNet-V2 [4] using a single input view.