

Light Field Neural Rendering

Mohammed Suhail^{1,2,*} Carlos Esteves⁴ Leonid Sigal^{1,2,3} Ameesh Makadia⁴

suhail33@cs.ubc.ca machc@google.com lsignal@cs.ubc.ca makadia@google.com

¹University of British Columbia ²Vector Institute for AI ³Canada CIFAR AI Chair ⁴Google

Abstract

Classical light field rendering for novel view synthesis can accurately reproduce view-dependent effects such as reflection, refraction, and translucency, but requires a dense view sampling of the scene. Methods based on geometric reconstruction need only sparse views, but cannot accurately model non-Lambertian effects. We introduce a model that combines the strengths and mitigates the limitations of these two directions. By operating on a four-dimensional representation of the light field, our model learns to represent view-dependent effects accurately. By enforcing geometric constraints during training and inference, the scene geometry is implicitly learned from a sparse set of views. Concretely, we introduce a two-stage transformer-based model that first aggregates features along epipolar lines, then aggregates features along reference views to produce the color of a target ray. Our model outperforms the state-of-the-art on multiple forward-facing and 360° datasets, with larger margins on scenes with severe view-dependent variations. Code and results can be found at [light-field-neural-rendering.github.io](https://github.com/light-field-neural-rendering).

1. Introduction

Synthesizing a novel view given a sparse set of images is a long-standing challenge in computer vision and graphics [10, 41, 42]. Recent advances in 3D neural rendering for view synthesis, in particular NeRF [31] and its successors [15, 16, 27, 33, 36, 58], have brought us tantalizingly close to the capability of creating photo-realistic images in complex environments. One reason for NeRF’s success is its implicit 5D scene representation which maps a 3D scene point and 2D viewing direction to opacity and color. In principle, such a representation could be perfectly suited to modeling view-dependent effects such as the non-Lambertian reflectance of specular and translucent surfaces. However, without regularization, this formulation permits degenerate solutions due to the inherent ambiguity between 3D sur-

* Work done while interning at Google.

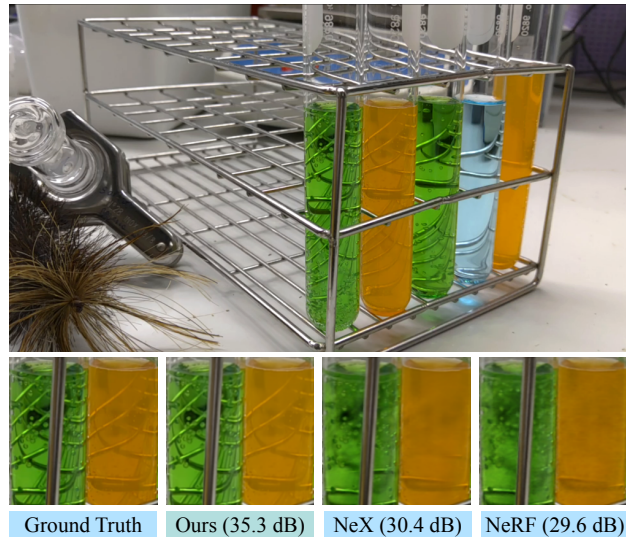


Figure 1. **Novel view synthesis.** On top is the target image to be rendered, from the *Lab* scene in the Shiny dataset [55]. Bottom row shows crops of novel views generated by our proposed model, NeX [55], and NeRF [31]. Unlike NeX and NeRF that fail to synthesize refractions on the test tube, our model almost perfectly reconstructs these complex view-dependent effects. We indicate the PSNR of the rendered images within parenthesis (higher is better). Images can be zoomed for detail.

face and radiance, where an incorrect shape (opacity) can be coupled with a high-frequency radiance function to minimize the optimization objective [61]. In practice, NeRF avoids such degenerate solutions through its neural architecture design, where the viewing direction is introduced only in the last layers of the MLP, thereby limiting the expressivity of the radiance function, which effectively translates to a smooth BRDF prior [61]. Thus, NeRF manages to avoid degenerate solutions at the expense of fidelity in non-Lambertian effects (Fig. 1 highlights this particular limitation of the NeRF model). Photo-realistic synthesis of non-Lambertian effects is one of the few remaining hurdles for neural rendering techniques.

In this paper, we formulate view synthesis as rendering a

sparingly observed light field. The 4D light field [26], which measures the radiance along rays in empty space, is often used for view synthesis [6, 24, 26]. Rendering a novel view from a densely sampled light field can be achieved with signal processing techniques (*e.g.*, interpolation) and without any model of the 3D geometry, but no such straightforward method exists with sparse light fields. From sparse images, rendering often utilizes additional 3D geometric constraints, such as predicted depth maps [24, 46], but performance is sensitive to accurate depth estimates which are difficult to obtain for non-Lambertian surfaces.

Motivated by these limitations, we introduce a novel method for rendering a sparse light field. Our neural rendering function operates in the style of image based rendering, where a target ray is synthesized using only observed rays from nearby views. In lieu of explicit 3D information, our transformer based rendering function is trained to fuse rays from nearby views exploiting an additional inductive bias in the form of a multi-view geometric constraint, namely the epipolar geometry. As shown in Fig. 1, our model is able to faithfully reconstruct the sharp details and lighting effect in the most challenging scene in the Shiny dataset [55].

Contributions. Our main contribution is the novel light field based neural view synthesis model, capable of photorealistic modeling of non-Lambertian effects (*e.g.*, specularities and translucency). To address the core challenge of sparsity of initial views, we leverage an inductive bias in the form of a multi-view geometric constraint, namely the epipolar geometry, and a transformer-based ray fusion. The resulting model produces higher fidelity renderings for forward-facing as well as 360° captures, compared to state-of-the-art, achieving up to 5 dB improvement in the most challenging scenes. Further, as a byproduct of our design, we can easily obtain dense correspondences and depth without further modifications, as well as transparent visualization of the rendering process itself. Through ablations we illustrate the importance of our individual design choices.

2. Related work

Light field rendering. Levoy and Hanrahan [26] defined the 4D light field as a function that specifies the radiance of any given ray in free space. They forewent geometric reasoning to directly synthesize novel views from input samples. Lumigraph rendering [18] exploits proxy geometry to counter aliasing effects that stem from irregularity or view under-sampling. Recent works [6, 24, 45, 46, 56] have explored learning based methods to light field rendering. These methods, however, either require dense input sampling [24], have limited range of motion [46] or are limited to simple scenes [45]. In this work, we focus on novel view synthesis for complex scenes with challenging non-Lambertian effects from a sparse set of viewpoints.

Neural scene representation. Representing shape and appearance of scenes using neural networks has recently gained immense popularity. *Explicit* representation-based methods use differentiable rendering to learn 3D representation such as point clouds [1, 40, 57], meshes [50] or voxels [28, 43] for the scene. *Implicit* representation-based methods represent scenes using continuous coordinate-based functions such as signed distance fields [2, 8, 17, 23, 59, 60] or occupancy fields [29, 34]. Scene Representation Networks [44] use a differentiable ray marching algorithm along with a continuous function that maps coordinates to features. NeRF [31] achieves photo-realistic rendering by learning a function that maps points along a ray to color and opacity followed by volumetric rendering. NeX [55] is a multiplane image-based scene representation that addresses NeRF’s difficulty to model large view dependent effects. However, NeX is still challenged in scenarios such as interference patterns caused by reflection or refraction through liquid. In this work, we introduce a model that can faithfully render novel views in the presence of complex view-based effects in scenarios where other methods fail (*e.g.*, Fig. 1). For a comprehensive survey of recent advances in neural rendering please refer to Tewari et al. [49].

Image-based rendering. Image-based rendering (IBR) methods [9, 12, 14, 54] are built on the notions that novel views can be rendered by “borrowing” pixel values from a given set of input images. Global geometry-based methods such as Hedman and Kopf [19], Hedman et al. [20], Riegler and Koltun [38, 39] rely on dense reconstruction from input views to obtain a global mesh for the scene. These meshes are used for projecting target rays onto nearby view for feature and color extraction. Other methods such as Chaurasia et al. [9], Penner and Zhang [35] infer depth maps using multi-view stereo methods to compute warping transforms from the given input to target viewpoints. Thies et al. [51] combine IBR with GAN-based image synthesis to learn view-dependent effects. To overcome difficulties caused by error in depth estimation, Choi et al. [11] estimate a depth uncertainty distribution to refine images.

Recently, Wang et al. [54] introduced IBRNet, a NeRF-based model that incorporates features from nearby views for rendering. Their architecture predicts colors for each point on the ray as weighted average of colors from neighboring views. The densities for each point are predicted by aggregating information from all other points using a single attention layer. Similarly for category-specific reconstruction, NerFormer [37] proposed to replace the MLP in NeRF-WCE [22] with a transformer model to allow for spatial reasoning. Our work crucially differs from these methods at their core: the rendering framework. Our method employs a light field representation, forgoing the need for volumetric rendering. Furthermore, we introduce a transformer-based model that first reasons about correspondences to ag-

gregate features along each epipolar line, then reasons about occlusion and lighting effects to aggregate features from multiple views to produce the final color.

3. Approach

Our goal is to synthesize novel views of a scene given a collection of input images available during both training and inference. Our design is guided by two key ideas, 1) using the four-dimensional parametrization of the light-field as input enables capturing view-dependent effects with high fidelity, and 2) enforcing constraints from multiple view geometry allows for view synthesis with sparse input views.

These ideas enable faithfully recovering illumination effects as in classical lightfield methods [18, 26], but require only a sparse view-sampling of the scene, as in geometry-based methods [12, 32] which traditionally struggle reproducing non-Lambertian effects. To implement them, we introduce an epipolar-geometric inductive bias in conjunction with a transformer-based architecture. Our model can render novel views from forward-facing photos as well as 360° scenes captured with cameras on a hemisphere.

In the following sections, we first introduce the light field representation, then an overview of the model, followed by a detailed description of the network architecture.

3.1. Light field parametrization

Light fields are functions on the space of oriented lines that associate a radiance value to a given ray. In free space, as the radiance along the ray remains constant, the space of rays has four degrees of freedom and can be parametrized by 4D vectors. We consider two distinct parametrizations of light field, the light slab [26] and the two-sphere [7].

Light slab. We adopt the light slab parametrization for forward-facing captures. A light slab consists of two parallel planes with their respective 2D coordinate systems (s, t) and (u, v) . Rays are then represented as a 4D tuple $r = (s, t, u, v)$ containing the coordinates of intersections with the two planes in their respective coordinate frames.

Two-sphere. For 360° scenes, we use the two-sphere parametrization [7] of the light field. Given a sphere bounding a scene, rays from the camera are represented using the colatitudes and longitudes at the two intersections with the sphere, $r = (\theta_1, \phi_1, \theta_2, \phi_2)$.

Given a four-dimensional light field ray parametrization r , we learn a neural rendering model f that maps the rays to radiance values.

To obtain the ray coordinates for a given pixel in homogeneous coordinates $x \in \mathbf{RP}^2$, from an image taken using a camera with intrinsics C and pose (extrinsics) $[R \ t]$, we first obtain the ray as a line ℓ in world coordinates parametrized by δ as $\ell(\delta) = -R^\top t + \delta R^\top C^{-1}x$, then solve for δ to obtain the intersections r with either the two planes or the

sphere. To render an image, we evaluate the model $f(r)$ for the rays associated to each target pixel.

3.2. Model overview

Optimizing a neural rendering model f that directly maps 4D light field coordinates to color fails to generalize to novel views when trained with a sparse set of input views (see Sec. 4.3 for quantitative evaluation).

To address this challenge, we introduce a model that incorporates a geometric inductive bias in the form of the epipolar constraints.

Given a target camera, we identify a set of neighboring views to be used to enforce multiple-view consistency. During training, this set is constructed by randomly choosing K views from a subset of N closest views. During inference, the closest K are chosen deterministically. We refer to the set of K chosen views as *reference views*.

Now given a target pixel x to be rendered, we obtain its ray parametrizations ℓ and r as described in Sec. 3.1, sample a sequence of P points $p_i = \ell(\delta_i)$ along the ray, and project each point to each reference view as $x_i^j = C_j[R_j \ t_j]p_i$, where $C_j, [R_j \ t_j]$ are the reference view camera intrinsics and extrinsics, respectively, and $1 \leq j \leq K$.

The collection $x^j = \{x_i^j\}_{1 \leq i \leq P}$ consists of points along the epipolar line of the target ray in the j^{th} reference view, and we refer to x_i^j as *epipolar points*. To each epipolar point, we associate its ray parametrization as described in Sec. 3.1, yielding the collections $r^j = \{r_i^j\}_{1 \leq i \leq P}$.

Epipolar feature aggregation. The first stage of our model, represented by the function f_1 , computes a feature representation per reference view by aggregating features associated to the epipolar points and target ray. We detail what those features are in the following sections. Conceptually, the first stage computes the set of features $\{z^j\}_{1 \leq j \leq K}$ where $z^j = f_1(r, r^j)$. This is loosely related to classical multiple view geometry, where we look for a correspondence to the target ray along the epipolar line. In our case, however, there is no visual representation of the target ray, so the model must learn to match the target ray coordinates with the available reference features, and the output is a feature vector representing the view j .

View feature aggregation. The second stage, represented by the function f_2 , predicts the target ray color by aggregating features associated to each reference view, given the target ray representation. Conceptually, the color for pixel x with associated ray r is predicted as $f(r) = f_2(r, \{f_1(r, r^j)\})$. This stage learns to reason about occlusion and illumination effects to combine information from all views and produce the target ray color.

3.3. Network architectures

One possible approach would be to model f_1 and f_2 as multi-layer perceptrons (MLP). However, this impedes the

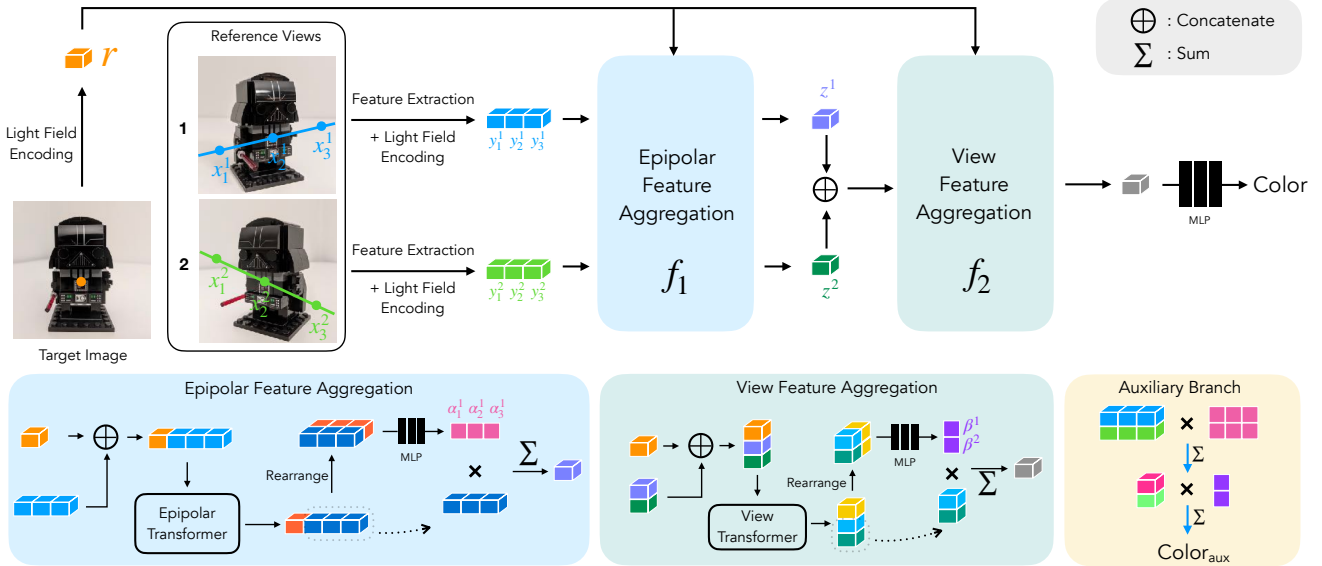


Figure 2. **Model Overview.** Given a target ray to render, we identify reference views and sample points along the epipolar lines corresponding to the target ray. Features of these epipolar points along with the light field coordinates of the target ray are inputs to the epipolar aggregation. This stage (blue), *independently* aggregates features along the epipolar lines for each reference view, producing reference view features. The reference view features along with the target ray are passed to the view aggregation stage (green), which combines the reference view features to predict the target ray color.

model from exploiting readily available relational information that can be extracted from the epipolar points, leading to sub-optimal performance (Sec. 4.3 quantifies this). Since inputs to f_1 are a *sequence* of epipolar points, and inputs to f_2 are a *set* of reference view features, we propose to use transformers, which excel in sequence and set modeling, to model both epipolar and view feature aggregation.

3.3.1 Epipolar feature transformer (f_1)

This transformer, highlighted in blue in Fig. 2, combines the features of points along the epipolar line based on the target rays.

The input is a sequence of $P + 1$ features, with P features from epipolar points and one from the target ray. The feature vector for the target ray is its own coordinates r . The feature for an epipolar point x_i^j is a concatenation of 1) ray coordinates r_i^j , 2) coordinates of p_i , the 3D point along r projected to x_i^j , 3) a learnable camera embedding k_j , 4) visual features v_i^j at x_i^j , obtained from a lightweight CNN, and 5) the color c_i^j at x_i^j .

Assuming the target pixel x matches to an epipolar point x_i^j , the corresponding point in the scene can be solved for and will have coordinates of p_i . Including it as an epipolar point feature also plays the role of positional encoding, since each point in the epipolar line correspond to some depth value along the query ray. This type of positional

encoding is richer than the typical 1D encoding used in sequence modeling [52], and more appropriate for modeling a 3D scene, as demonstrated in Sec. 4.3.

We further apply Fourier features [31, 48] positional encoding to facilitate learning of high-frequency functions. This operation is performed by γ_r for ray coordinates and γ_p for point coordinates, see Sec. 4.1 for details. To summarize, each epipolar point x_i^j is represented by a feature

$$y_i^j = [\gamma_r(r_i^j) \parallel \gamma_p(p_i) \parallel k_j \parallel v_i^j \parallel c_i^j], \quad (1)$$

where \parallel denotes concatenation. The epipolar transformer for view j will take as inputs $[\gamma_r(r), \{y_i^j\}_{1 \leq i \leq P}]$. A linear layer first projects the features to the same dimension, then a self-attention transformer is applied to the whole sequence.

We aggregate the P outputs corresponding to the epipolar points (\tilde{y}_i^j) to obtain the *reference view features*. The aggregation is a weighted average, with the weights computed using an attention mechanism similar to the Graph Attention Networks (GAT) [53] as follows,

$$\alpha_i^j = \frac{\exp(W_1 [\tilde{r} \parallel \tilde{y}_i^j])}{\sum_k \exp(W_1 [\tilde{r} \parallel \tilde{y}_k^j])}, \quad (2)$$

where \tilde{r} are the output features of the target ray, and W_1 are learned weights. The first stage is completed by repeating $z^j = f_1(r, r^j) = \sum_{i=1}^P \alpha_i^j \tilde{y}_i^j$ for all views $1 \leq j \leq K$.

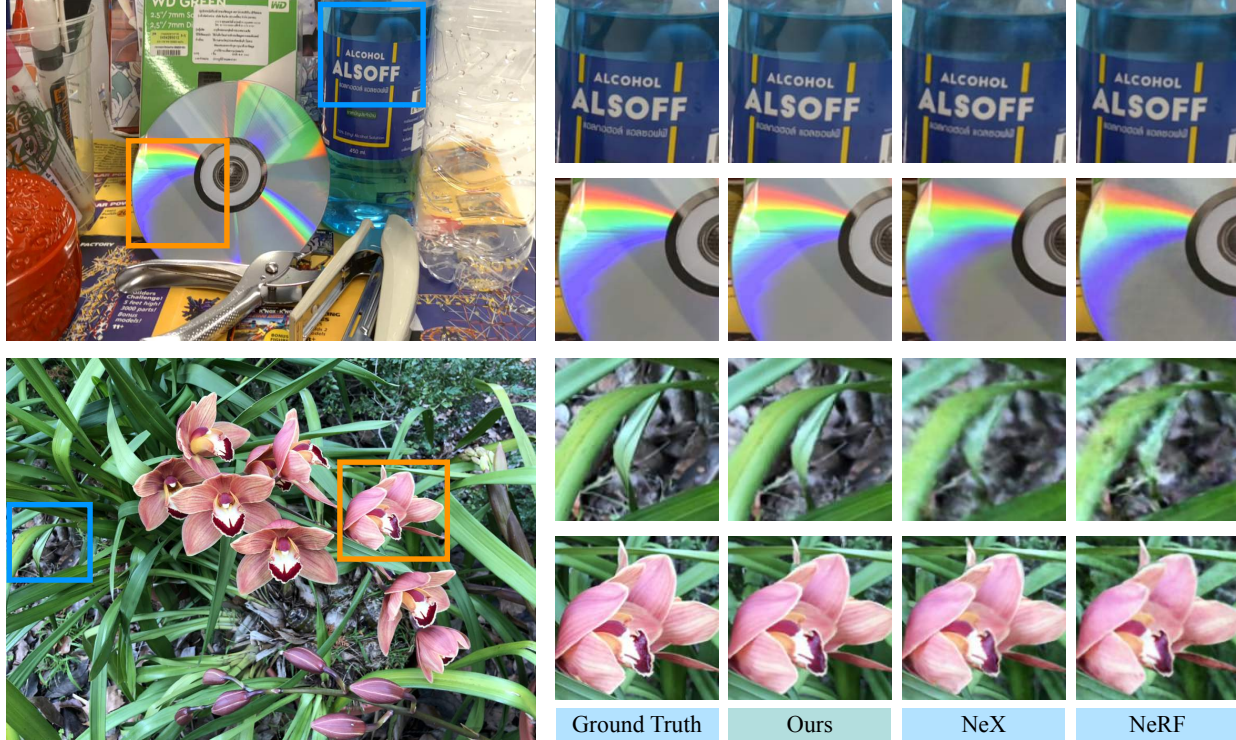


Figure 3. **Qualitative Comparison.** Top: results on the *CD* scene from Shiny dataset [55]. Our method is able to retrieve sharper details in the reflections on the bottle (e.g., the top-left of the insets) as well as the interference patterns on the compact disk (e.g., rainbow and reflection on the top-right). Bottom: results on the *Orchids* scene from the real-forward facing (RFF) dataset [30]. Our method recovers more accurately the shape of the leaves. We also observe sharper texture on the leaves as well as the petals.

3.3.2 View feature transformer (f_2)

This transformer, highlighted in green in Fig. 2, takes the target ray and the set of features for each reference view. The input sequence is now $[\gamma_r(r), \{z^j\}_{1 \leq j \leq K}]$, where z^j are the reference view features computed by the first stage, and the output is a single feature vector for the target ray. We use the same self-attention transformer architecture as the epipolar feature aggregator. The transformer output sequence $[\hat{r}, \{\tilde{z}^j\}_{1 \leq j \leq K}]$ is aggregated with a weighted average using the same idea as the previous section. We compute the weights β^j with learnable weights W_2 ,

$$\beta^j = \frac{\exp(W_2[\hat{r} \parallel \tilde{z}^j])}{\sum_k \exp(W_2[\hat{r} \parallel \tilde{z}^k])}, \quad (3)$$

then the output of this stage is the target ray feature $\sum_{k=1}^K \beta^k \tilde{z}^k$, which is linearly projected and passed by a sigmoid to produce the pixel color prediction c .

3.4. Loss

During training, we minimize the \mathcal{L}_2 loss between the observed and predicted colors. We additionally include an auxiliary loss to encourage the attention weights for the

epipolar points (α_i^j) and reference views (β^j) to be interpretable, in the sense that high values of α_i^j suggest a valid match to the target ray, while low values of β^j might indicate occlusion. This auxiliary loss also leads to more accurate renderings (see Sec. 4.3). To compute it, we use the attention weights to combine reference pixel colors and make a second color prediction as

$$c_{\text{aux}} = \sum_j \beta^j \left(\sum_i \alpha_i^j c_i^j \right), \quad (4)$$

where c_i^j is the color of the epipolar point x_i^j . The auxiliary loss is then defined as the \mathcal{L}_2 loss between c_{aux} and the ground truth. The effect of this loss is two-fold: 1) it incentivizes weights α_i^j to have lower entropy to avoid blurry predictions in the auxiliary branch, and 2) it encourages weights β^j to be high for unoccluded views.

4. Experiments

We show quantitative and qualitative comparisons against state-of-the-art methods for novel view synthesis. We also perform an ablation study to analyze the effectiveness of the components introduced in our method.

Model	PSNR [dB] \uparrow	SSIM \uparrow	LPIPS \downarrow	Avg. \downarrow
LLFF [30]	24.41	0.863	0.211	0.0656
NeRF [31]	26.76	0.883	0.246	0.0562
IBRNet [54]	26.73	0.851	0.175	0.0523
NeX [55]	27.26	0.904	0.178	0.0473
Ours	28.26	0.920	0.062	0.0297

Table 1. Results for the real forward-facing (RFF) dataset [30].

4.1. Implementation details

Network architecture. We use similar transformer architectures as the ones recently introduced for vision related tasks [13]. Each block consists of a single-headed self-attention layer and an MLP with Gaussian error linear unit (GELU) activation [21]. A residual connection is applied at every block, followed by a LayerNorm (LN) [3]. Each transformer has 8 blocks and the internal feature size is 256. The visual features v_i^j are produced by a single convolutional layer with 5×5 filters and 32 channels.

Positional encoding. Following prior work [31, 48], we use Fourier features to encode input coordinates to facilitate learning the high-frequency components required for accurate rendering. For the light slab parametrization and the 3D points p_i , we positionally-encode each ray coordinate [31] as $\gamma_r(w) = \gamma_p(w) = \{\sin(2^k w)\} \cup \{\cos(2^k w)\}$ for $0 \leq k \leq 4$. For the two-sphere parametrization, we found it beneficial to use a positional encoding based on evaluating the spherical harmonics at the points (θ_1, ϕ_1) and (θ_2, ϕ_2) , see the appendix for details. The learnable camera embeddings k_j are 256-dimensional.

Training/inference details. In each training step, we randomly choose a target image and sample a batch of random rays from it. The batch sizes are 4096 for the forward-facing datasets and 8192 for Blender. We train for 250 000 iterations with the Adam optimizer [25] and a linear learning rate decay schedule with 5000 warm-up steps. For inference, we sample contiguous blocks of rays to make a batch. Training on a Blender scene takes around 23 hours on a 32-core TPUv3 slice. Rendering an 800×800 image then takes around 9.2 seconds.

4.2. Results

We compare our method with LLFF [30], NeRF [31], IBRNet [54], NeX [55] and Mip-NeRF [4]. We compare against Mip-NeRF only on the Blender dataset because for forward-facing captures, as noted by Barron et al. [4, Appx D], Mip-NeRF performs on par with NeRF.

Metrics. To measure the performance of our model we use three widely adopted metrics: peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM), and

Model	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Avg. \downarrow
NeRF [31]	25.60	0.851	0.259	0.0651
NeX [55]	26.45	0.890	0.165	0.0499
IBRNet [†] [54]	26.50	0.863	0.122	0.0468
Ours	27.34	0.907	0.045	0.0294

Table 2. Results for the Shiny dataset from NeX [55].

Model	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Avg. \downarrow
NeRF [31]	31.01	0.953	0.050	0.0194
IBRNet [54]	28.14	0.942	0.072	0.0299
Mip-NeRF [4]	33.09	0.961	0.043	0.0161
Ours	33.85	0.981	0.024	0.0110

Table 3. Results for the Blender dataset from NeRF [31].

the learned perceptual image patch similarity (LPIPS) [62]. Following [4], we additionally report the geometric mean of $10^{-\text{PSNR}/10}$, $\sqrt{1 - \text{SSIM}}$ and LPIPS, which provides a summary of three metrics for easier comparison. We report the averages of each metric over all the scenes in each dataset. Please refer to the appendix for a scene-wise breakdown of the results.

4.2.1 Real-forward-facing (RFF) dataset

The RFF dataset introduced by Mildenhall et al. [30] consists of 8 forward facing captures of real-world scenes using a smartphone. For our experiments, we use the same resolution and train/test splits as NeRF [31].

Table 1 reports the average metrics across all 8 scenes in the RFF dataset. We show qualitative comparisons on the *Orchids* scene in Fig. 3. Compared to the baselines, our method retains sharper detailed textures and produces consistent shape boundaries on the leaves and petals.

4.2.2 Shiny dataset

The RFF dataset mostly consists of diffuse scenes with little view-dependent effects. The Shiny dataset introduced in NeX [55] presents 8 scenes with challenging view-dependent effects, captured by forward-facing cameras. We use the same image resolution and splits as NeX.

We compare our model against NeX, IBRNet and NeRF on the Shiny dataset in Tab. 2. We report the average scores across all scenes in the dataset. Our model consistently improves over the state-of-the-art in all metrics. We show qualitative analysis of rendering on a test view from the *CD* scene in Fig. 3. Our model is able to reconstruct the interference patterns on the disk and reflections on the bottle with

Model	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Avg. \downarrow
Vanilla-NLF	17.39	0.614	0.516	0.1802
1-MLP	21.33	0.774	0.208	0.0900
2-MLP	26.16	0.896	0.076	0.0390
No CNN (v_i^j)	27.43	0.910	0.057	0.0314
No 3D Coordinates	28.17	0.920	0.047	0.0273
No LCE (k_j)	28.23	0.926	0.045	0.0264
Mean Pooling	28.39	0.929	0.043	0.0255
No Auxiliary Loss	28.43	0.931	0.043	0.0253
Ours	28.78	0.934	0.038	0.0235

Table 4. Ablation study on the RFF dataset [30], with 25% of the original resolution (504 × 378). Refer to Sec. 4.3 for details.

higher level of detail as compared to baselines.

4.2.3 Blender dataset

Our model is capable of rendering novel views of 360° scenes. To evaluate this case, we use the synthetic dataset introduced by Mildenhall et al. [31]. Each scene consists of 800 × 800 resolution images rendered from viewpoints randomly sampled on a hemisphere around the object.

Table 3 reports the average performance across all scenes in the Blender dataset. Our model improves over NeRF, IBRNet and Mip-NeRF on all metric and achieves new state-of-the-art results. On the materials scene, which contains reflections on metallic balls, we observe an improvement of around 4 dB on the PSNR metric when compared to Mip-NeRF. We present the full table along with qualitative comparisons in the appendix.

4.3. Ablation studies

To validate the effectiveness of different design decisions, we run the following ablation experiments.

Geometric inductive bias. We train a model, called ‘Vanilla-NLF’, that uses an MLP to predict the color of a ray given only its light field representation, without consideration of the scene geometry in form of epipolar constraints.

Transformers vs MLPs. We train variations of our model replacing the transformers with MLPs. In the first variant, we replace each of the epipolar and view transformers with MLPs (‘2-MLP’). The second variant replaces both transformers by a single MLP that takes as input all epipolar point features from all reference views along with the target ray and directly predicts the color (‘1-MLP’). We detail the architectures in the appendix. We run a sweep over the number of layers for these MLPs and report performance of the best model.

† We fine tune the pretrained model available at <https://github.com/googleinterns/IBRNet> on each scene in Shiny.

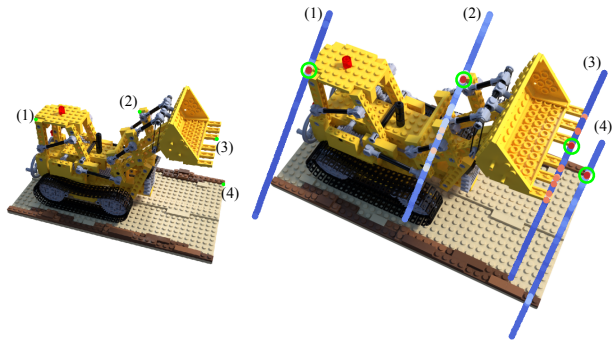


Figure 4. **Correspondence Distribution.** The per-point attention weights learned by our model indicate potential correspondences to the target ray. We visualize four target rays and one reference view. The weights are in log-scale, from blue to red. The green circles highlight the point with highest correspondence probability.

Model Components. To probe the efficacy of the different components, we train ablated models 1) without visual features v_i^j (‘No CNN’), 2) without 3D coordinates p_i , 3) without the learnable camera embedding k_j (‘No LCE’), 4) replacing the attention based aggregation with mean pooling, and 5) removing the auxiliary loss term.

Table 4 reports the ablation results. All models are trained on images from the RFF datasets downsampled to 25% of the original resolution (504 × 378). We use the average metrics across all the scenes for comparison.

4.4. Interpreting the model

The use of transformers and epipolar geometry in our model permits interpretation of the results via the attention weights. We demonstrate this by extracting correspondences and depth maps. Also, our use of a four-dimensional light field representation enables the construction of epipolar-plane images (EPI) [5], which are interpretable reconstructions of the scene geometry.

Dense correspondence. We can extract potential correspondences between a target ray and a reference view j by finding the largest attention weights in $\{\alpha_i^j\}_{1 \leq i \leq P}$. Figure 4 shows the weight distribution of putative correspondences over the epipolar line for four points of interest. For points (1) and (4) we observe unimodal distributions with peaks at the point of correspondence. For point (2) we notice some uncertainty, while for point (3), the distribution is multi-modal with peaks around each blade with the highest peak near the correct correspondence.

Disparity map. Since each epipolar point corresponds to the projection of the target ray at a certain depth (the putative depth), we can use the correspondence distributions to estimate a depth map for a target ray. We first extract all the epipolar point (α_i^j) and reference view (β^j) attention weights as described in Sec. 3.3, then compute a weighted

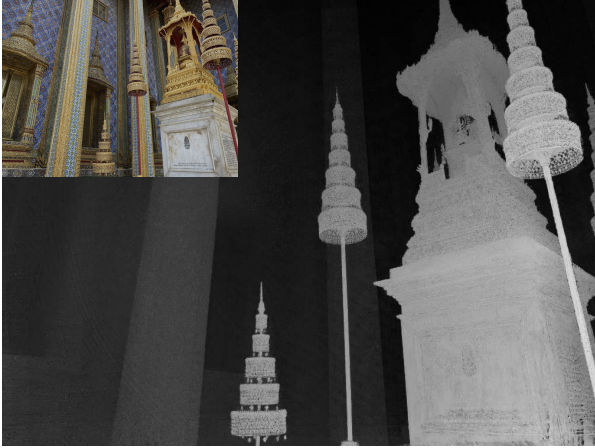


Figure 5. **Disparity Map.** The per-point and per-view attention weights learned by our model can be applied to estimate a disparity map by aggregating the putative depths of each epipolar point on each reference view, for each target ray.

average of putative depths, equivalent to applying Eq. (4) with colors replaced by depths. Figure 5 shows an example of a the disparity map obtained for a test view of the *Crest* scene from the Shiny dataset [55].

Epipolar-plane images (EPI). For a 4D light slab representation, we construct the EPI by querying our model with two fixed and two variable coordinates, resulting in a 2D color image. Physically, this corresponds to moving the camera along a 1D trajectory, stacking the images of a line segment parallel to the trajectory. EPIs encode information about specularities and scene geometry, where diffuse points appear as lines and specular points appear as curves. We show the epipolar slices for the *CD* and *Flower* scene in Fig. 6. The *Flower* scene is predominantly diffuse so we observe lines of varying slopes in the EPI, with slopes inversely proportional to the depth. For the *CD* scene, in addition to lines, we observe curves at regions corresponding to the interference pattern on the disk. This is due to the change in virtual apparent depth of the specular points with change in view point [30, 47].

5. Limitations

Since our method relies on implicitly finding the correspondences for a target pixel in nearby views, it is challenged by texture-less thin repeating structures. As shown in Fig. 7, our model produces fuzzy details on the grill-like structure in the *Tools* scene from Shiny and on the wire-mesh on the microphone from Blender.

Transformers are computationally expensive, resulting in slow training and inference times. Our method is around 8 times slower than Mip-NeRF [4] on same hardware. Our model does, however, compare favorably in terms of speed



Figure 6. **Epipolar-plane images (EPI).** Our model represents the 4D light field, so constructing EPIs is natural. Each EPI vertically stacks images along the blue line, while the camera moves parallel to the blue line. Different depths show as lines of different slopes in the EPI, while view-dependent effects show as curves.

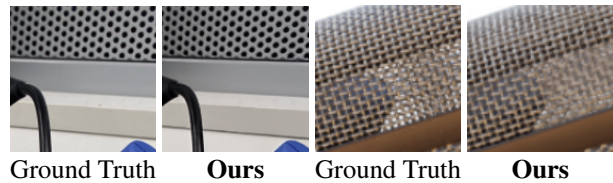


Figure 7. **Failure Cases.** Our model is challenged by texture-less thin repeating structures. Left: It produces distorted circles on the grill-like structure in the *Tools* scene from Shiny [55]. Right: Similar distortions appear for the *Mic* scene from Blender [31].

against other transformer-based models. For example, NerFormer [37] takes around 180 s to render an 800×800 image on a single V100 GPU, while our method takes 60 s to 70 s on the same hardware. We notice that our model suffers from overhead of 1) random memory access on device and 2) data transfer between host and device. We believe that it can be made more efficient with some engineering effort.

6. Conclusion

We present a light field based neural rendering method for novel-view synthesis. Unlike prior volumetric rendering methods, our proposed model can naturally handle real-world illumination effects by learning the light field over a four-dimensional space. To address the dense sampling dependency of light field rendering, we introduced a two-stage framework that incorporates geometric inductive bias in the form of epipolar constraints. Our model leads to significant improvement over previous state-of-the-art model for view synthesis especially for scenes with challenging view-dependent effects. Finally, the design of our model allows extracting dense correspondences, disparity maps and epipolar-plane images without any additional training.

References

- [1] K.-A. Aliev, A. Sevastopolsky, M. Kolos, D. Ulyanov, and V. Lempitsky. Neural point-based graphics. In *European Conference on Computer Vision (ECCV)*, pages 696–712, 2020. [2](#)
- [2] M. Atzmon, N. Haim, L. Yariv, O. Israelov, H. Maron, and Y. Lipman. Controlling neural level sets. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. [2](#)
- [3] L. J. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *CoRR*, abs/1607.06450, 2016. [6](#)
- [4] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. [6](#), [8](#)
- [5] R. C. Bolles, H. H. Baker, and D. H. Marimont. Epipolar-plane image analysis: An approach to determining structure from motion. *International Journal of Computer Vision (IJCV)*, 1(1):7–55, 1987. [7](#)
- [6] C. Buehler, M. Bosse, L. McMillan, S. Gortler, and M. Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer Graphics and Interactive Techniques*, pages 425–432, 2001. [2](#)
- [7] E. Camahort, A. Leros, and D. Fussell. Uniformly sampled light fields. In *Eurographics Workshop on Rendering Techniques*, pages 117–130, 1998. [3](#)
- [8] R. Chabra, J. E. Lenssen, E. Ilg, T. Schmidt, J. Straub, S. Lovegrove, and R. Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *European Conference on Computer Vision (ECCV)*, pages 608–625, 2020. [2](#)
- [9] G. Chaurasia, S. Duchene, O. Sorkine-Hornung, and G. Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM Transactions on Graphics (TOG)*, 32(3):1–12, 2013. [2](#)
- [10] S. E. Chen and L. Williams. View interpolation for image synthesis. In *Proceedings of the 20th annual conference on Computer Graphics and Interactive Techniques*, pages 279–288, 1993. [1](#)
- [11] I. Choi, O. Gallo, A. Troccoli, M. H. Kim, and J. Kautz. Extreme view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7781–7790, 2019. [2](#)
- [12] P. E. Debevec, C. J. Taylor, and J. Malik. Modeling and rendering architecture from photographs: A hybrid geometry- and image-based approach. In *Proceedings of the 23rd annual conference on Computer Graphics and Interactive Techniques*, pages 11–20, 1996. [2](#), [3](#)
- [13] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations (ICLR)*, 2021. [6](#)
- [14] R. Du, M. Chuang, W. Chang, H. Hoppe, and A. Varshney. Montage4d: Interactive seamless fusion of multiview video textures. *ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games*, 2018. [2](#)
- [15] Y. Du, Y. Zhang, H.-X. Yu, J. B. Tenenbaum, and J. Wu. Neural radiance flow for 4d view synthesis and video processing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14324–14334, 2021. [1](#)
- [16] G. Gafni, J. Thies, M. Zollhofer, and M. Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8649–8658, 2021. [1](#)
- [17] K. Genova, F. Cole, A. Sud, A. Sarna, and T. Funkhouser. Local deep implicit functions for 3d shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4857–4866, 2020. [2](#)
- [18] S. J. Gortler, R. Grzeszczuk, R. Szeliski, and M. F. Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 43–54, 1996. [2](#), [3](#)
- [19] P. Hedman and J. Kopf. Instant 3d photography. *ACM Transactions on Graphics (TOG)*, 37(4):1–12, 2018. [2](#)
- [20] P. Hedman, S. Alsisan, R. Szeliski, and J. Kopf. Casual 3d photography. *ACM Transactions on Graphics (TOG)*, 36(6):1–15, 2017. [2](#)
- [21] D. Hendrycks and K. Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. [6](#)
- [22] P. Henzler, J. Reizenstein, P. Labatut, R. Shapovalov, T. Ritschel, A. Vedaldi, and D. Novotny. Unsupervised learning of 3d object categories from videos in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4700–4709, 2021. [2](#)
- [23] C. Jiang, A. Sud, A. Makadia, J. Huang, M. Nießner, T. Funkhouser, et al. Local implicit grid representations for 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6001–6010, 2020. [2](#)
- [24] N. K. Kalantari, T.-C. Wang, and R. Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Transactions on Graphics (TOG)*, 35(6):1–10, 2016. [2](#)
- [25] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, 2015. [6](#)
- [26] M. Levoy and P. Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer Graphics and Interactive Techniques*, pages 31–42, 1996. [2](#), [3](#)
- [27] D. B. Lindell, J. N. Martel, and G. Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14556–14565, 2021. [1](#)
- [28] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Transactions on Graphics (TOG)*, 38(4):65:1–65:14, 2019. [2](#)
- [29] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,

- pages 4460–4470, 2019. [2](#)
- [30] B. Mildenhall, P. P. Srinivasan, R. Ortiz-Cayon, N. K. Kalantari, R. Ramamoorthi, R. Ng, and A. Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4): 1–14, 2019. [5](#), [6](#), [7](#), [8](#)
- [31] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*, pages 405–421, 2020. [1](#), [2](#), [4](#), [6](#), [7](#), [8](#)
- [32] P. Narayanan, P. W. Rander, and T. Kanade. Constructing virtual worlds using dense stereo. In *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*, pages 3–10. IEEE, 1998. [3](#)
- [33] M. Oechsle, S. Peng, and A. Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021. [1](#)
- [34] S. Peng, M. Niemeyer, L. Mescheder, M. Pollefeys, and A. Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision (ECCV)*, pages 523–540, 2020. [2](#)
- [35] E. Penner and L. Zhang. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)*, 36(6):1–11, 2017. [2](#)
- [36] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10318–10327, 2021. [1](#)
- [37] J. Reizenstein, R. Shapovalov, P. grid, L. Sbordone, P. Labatut, and D. Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (CVPR)*, pages 10901–10911, 2021. [2](#), [8](#)
- [38] G. Riegler and V. Koltun. Free view synthesis. In *European Conference on Computer Vision (ECCV)*, pages 623–640, 2020. [2](#)
- [39] G. Riegler and V. Koltun. Stable view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12216–12225, 2021. [2](#)
- [40] D. Rückert, L. Franke, and M. Stamminger. Adop: Approximate differentiable one-pixel point rendering. *arXiv preprint arXiv:2110.06635*, 2021. [2](#)
- [41] J. Shade, S. Gortler, L.-w. He, and R. Szeliski. Layered depth images. In *Proceedings of the 25th annual conference on Computer Graphics and Interactive Techniques*, pages 231–242, 1998. [1](#)
- [42] H. Shum and S. B. Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing*, volume 4067, pages 2–13, 2000. [1](#)
- [43] V. Sitzmann, J. Thies, F. Heide, M. Nießner, G. Wetzstein, and M. Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2437–2446, 2019. [2](#)
- [44] V. Sitzmann, M. Zollhöfer, and G. Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems*, 2019. [2](#)
- [45] V. Sitzmann, S. Rezkikov, W. T. Freeman, J. B. Tenenbaum, and F. Durand. Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. [2](#)
- [46] P. P. Srinivasan, T. Wang, A. Sreelal, R. Ramamoorthi, and R. Ng. Learning to synthesize a 4d rgbd light field from a single image. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2243–2251, 2017. [2](#)
- [47] R. Swaminathan, S. B. Kang, R. Szeliski, A. Criminisi, and S. K. Nayar. On the motion and appearance of specularities in image sequences. In *European Conference on Computer Vision (ECCV)*, pages 508–523, 2002. [8](#)
- [48] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. [4](#), [6](#)
- [49] A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, Y. Wang, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi, et al. Advances in neural rendering. *arXiv preprint arXiv:2111.05849*, 2021. [2](#)
- [50] J. Thies, M. Zollhöfer, and M. Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. [2](#)
- [51] J. Thies, M. Zollhöfer, C. Theobalt, M. Stamminger, and M. Nießner. IGNOR: Image-guided neural object rendering. *International Conference on Learning Representations (ICLR)*, 2020. [2](#)
- [52] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [4](#)
- [53] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph Attention Networks. *International Conference on Learning Representations (ICLR)*, 2018. [4](#)
- [54] Q. Wang, Z. Wang, K. Genova, P. P. Srinivasan, H. Zhou, J. T. Barron, R. Martin-Brualla, N. Snavely, and T. Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4690–4699, 2021. [2](#), [6](#)
- [55] S. Wizadwongsa, P. Phongthawee, J. Yenphraphai, and S. Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8534–8543, 2021. [1](#), [2](#), [5](#), [6](#), [8](#)
- [56] G. Wu, Y. Liu, L. Fang, and T. Chai. Revisiting light field rendering with deep anti-aliasing neural network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [2](#)
- [57] M. Wu, Y. Wang, Q. Hu, and J. Yu. Multi-view neural human rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages

1682–1691, 2020. 2

- [58] W. Xian, J.-B. Huang, J. Kopf, and C. Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9421–9431, 2021. 1
- [59] Q. Xu, W. Wang, D. Ceylan, R. Mech, and U. Neumann. DISN: Deep implicit surface network for high-quality single-view 3d reconstruction. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [60] L. Yariv, J. Gu, Y. Kasten, and Y. Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [61] K. Zhang, G. Riegler, N. Snavely, and V. Koltun. NeRF++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 1
- [62] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 6