# Learning Neural Light Fields with Ray-Space Embedding

Benjamin Attal*
Carnegie Mellon University

Jia-Bin Huang
Meta

Michael Zollhöfer
Reality Labs Research

Johannes Kopf
Meta

Changil Kim
Meta

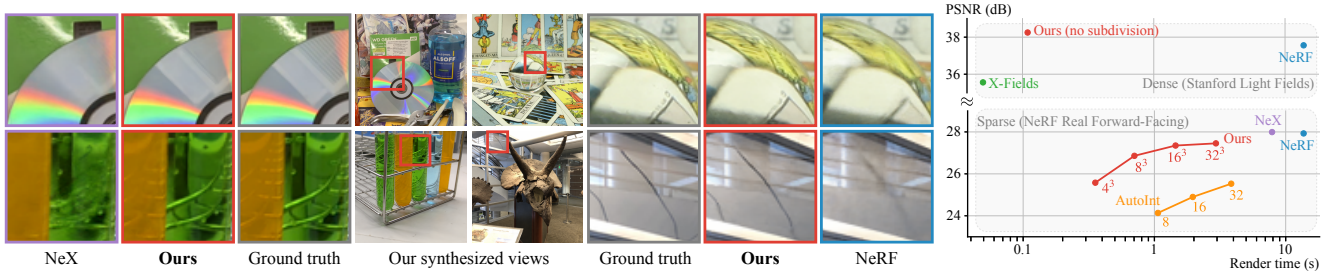https://neural-light-fields.github.io



Figure 1. We present *Neural Light Fields with Ray-Space Embedding*, which map rays directly to integrated radiance. Our representation allows for high-quality view synthesis with faithful reconstruction of complex view dependence (*left*). We are able to synthesize novel views in a fraction of the time required for the current state-of-the-art approaches, with comparable memory footprint. Our method achieves a more favorable trade-off between quality, speed, and memory than existing approaches for both dense and sparse light fields (*right*). See more visual comparisons including video results in our *project website* linked above.

## Abstract

*Neural radiance fields (NeRFs) produce state-of-the-art view synthesis results, but are slow to render, requiring hundreds of network evaluations per pixel to approximate a volume rendering integral. Baking NeRFs into explicit data structures enables efficient rendering, but results in large memory footprints and, in some cases, quality reduction. Additionally, volumetric representations for view synthesis often struggle to represent challenging view dependent effects such as distorted reflections and refractions. We present a novel neural light field representation that, in contrast to prior work, is fast, memory efficient, and excels at modeling complicated view dependence. Our method supports rendering with a single network evaluation per pixel for small baseline light fields and with only a few evaluations per pixel for light fields with larger baselines. At the core of our approach is a ray-space embedding network that maps 4D ray-space into an intermediate, interpolable latent space. Our method achieves state-of-the-art quality on dense forward-facing datasets such as the Stanford Light Field dataset. In addition, for forward-facing scenes with sparser inputs we achieve results that are competitive with NeRF-based approaches while providing a better speed/quality/memory trade-off with far fewer network evaluations.*

## 1. Introduction

View synthesis is an important problem in computer vision and graphics. Its goal is to photorealistically render a scene from unobserved camera poses, given a few posed input images. Existing approaches solve this problem by optimizing some underlying representation of the scene's appearance and geometry and then rendering this representation from novel views.

View synthesis has recently experienced a renaissance with an explosion of interest in *neural scene representations* — see Tewari et al. [54,55] for a snapshot of the field. *Neural radiance fields* [31] are perhaps the most popular of these neural representations, and methods utilizing them have recently set the state-of-the-art in rendering quality for view synthesis. A radiance field is a 5D function that maps a 3D point $\mathbf{x}$ and 3D direction $\vec{\omega}$ (with only 2 degrees of freedom) to the radiance leaving $\mathbf{x}$ in direction $\vec{\omega}$, as well as the density of the volume at point $\mathbf{x}$. A neural radiance field or NeRF represents this function with a neural network. Because volume rendering a NeRF is differentiable, it is straightforward to optimize by minimizing the difference between ground truth views at known camera poses and their reconstructions.

The main drawback of neural radiance fields is that volume rendering requires many samples and thus many neural network evaluations per ray to approximate a volume rendering accurately. Thus, rendering from a NeRF is usually quite slow. Various approaches exist for baking or caching neural

---

* This work was done while Benjamin was an intern at Meta.

radiance fields into explicit data structures to improve efficiency [12, 16, 41, 68]. Some approaches, concurrent to this work, learn color and density directly on a voxel grid, which improves both training and rendering speed [32, 52, 67]. However, the storage cost for explicit representations is much higher than a NeRF. Further, the baking procedure itself sometimes leads to a loss in resulting view synthesis quality for baking methods. Other methods aim to reduce the number of neural network evaluations per ray by representing radiance only on surfaces [18, 33]. These methods predict new images with only a few evaluations per ray. However, their quality is contingent on either ground truth geometry or reasonable geometry estimates, which are not always available.

A *light field* [13, 21] is the integral of a radiance field. It maps ray parameters directly to the integrated radiance along that ray. Thus, only one look-up of the underlying representation is required to determine the color of a ray; hence one evaluation per pixel, unlike hundreds of evaluations required by a NeRF. A common assumption for light fields is that this integral remains the same no matter the ray origin (i.e., radiance is constant along rays), which holds when the convex hull of the scene geometry does not contain any viewpoints used for rendering [21]. Given this assumption, a light field is a function of a ray on a 4D ray space.

In this paper, we show how to learn *neural light fields*. Since coordinate-based neural representations have been successfully employed to learn radiance fields from a set of ground truth images, one might think that they could be useful for representing and learning light fields as well. However, we show that learning light fields is significantly more challenging than learning radiance fields. Using the same neural network architecture as in NeRF to parameterize a light field leads to poor interpolation quality for view synthesis.

While a radiance field is a 5D function, an essential ingredient of NeRF is that it learns the scene geometry as a density field in 3D space. Additionally, a NeRF's learned appearance is closer to a 3D function than a 5D function since the network is late-conditioned on viewing directions [69]. This makes NeRFs easy to optimize but also means that they can struggle to represent complex view-dependent effects such as reflections and refractions [61] which violate multi-view color constraints.

On the other hand, we face the problem of learning a function defined on a 4D ray space from only partial observations — input training images only cover a few 2D slices of the full 4D space. At the same time, NeRF enjoys multiple observations of most 3D points. Further, light fields do not entail any form of scene geometry, which allows them to capture complex view dependence but poses significant challenges in interpolating unseen rays in a geometrically meaningful way. Existing methods address these issues by sacrificing

view interpolation [9], leveraging data driven priors [49], or relying on strong supervision signals [23].

In order to deal with these challenges, we employ a novel ray-space embedding network that re-maps the input ray-space into an embedded latent space. This facilitates both the registration of rays observing same 3D points and the interpolation of unobserved rays, which leads to better memorization and view synthesis at the same time. The embedding network alone already provides state-of-the-art view synthesis quality for densely sampled inputs (such as the Stanford light fields [60]). However, it does not interpolate well in sparser input sequences (such as those from Real Forward-Facing [31]). Thus, we represent such scenes with a set of local light fields, where each local light field is less prone to large depth and texture changes. Each local light field has to learn a simpler embedding at the price of several more network evaluations per ray.

We evaluate our method for learning neural light fields in sparse and dense regimes, with and without subdivisions. We compare to state-of-the-art view-synthesis methods and show that our approach achieves comparable or better view synthesis quality in both regimes, in a fraction of render time that other methods require, while still maintaining their small memory footprint (Figure 1).

In summary, our contributions are:

1. A novel neural light field representation that employs a ray-space embedding network and achieves state-of-the-art quality for small-baseline view synthesis without any geometric constraints.

2. A subdivided neural light field representation for large baseline light fields that leads to a good trade-off in terms of the number of network queries vs. quality, which can be optimized to achieve comparable performance to NeRF [31] and NeX [61] for sparse real-world scenes.

3. Improved capture of view-dependent appearance in both sparse and dense regimes (e.g., complicated reflections and refractions) that existing volume-based methods such as NeRF [31] and NeX [61] struggle to represent.

## 2. Related Work

**Light Fields.** Light fields [21] or Lumigraphs [13] represent the radiance entering all positions in 3D space, in all viewing directions. By resampling the rays from the discrete set of captured images, one can query arbitrary rays in a light field and synthesize photorealistic novel views of a static scene without necessitating geometric reconstruction. While light fields have the flexibility to represent complicated, non-Lambertian appearance, this often requires high sampling rates, and thus capturing and storing an excessive number of images. Much effort has been devoted to extending light fields to sparse, unstructured input sequences [4, 7] and synthesizing (i.e. interpolating) in-between views for

such sequences [17, 45, 59, 62, 63, 66]. Very recently, neural implicit representations have been applied to modeling light fields for better memory compactness [9] and faster rendering for view synthesis [23, 49]. Like these methods, our work leverages neural light fields as its core representation for view synthesis. Unlike prior work, which sacrifices view interpolation [9] or leverages meta-learning for strong data-driven priors [49], we enable view synthesis without explicit priors, using per-scene training only.

**Representations for View Synthesis.** Various 3D representations have been developed for rendering novel views from a sparse set of captured images. Image-based rendering techniques leverage proxy geometry to warp and blend source image content into novel views [4, 6, 38, 47]. Recent advances in image-based rendering include learning-based disparity estimation [10, 17], blending [15, 42], and image synthesis via 2D CNNs [42, 43]. Other commonly used scene representations include voxels [26, 50], 3D point clouds [1, 19, 44] or camera-centric layered 3D representations, such as multi-plane images [3, 30, 57, 61, 71] or layered depth images [14, 20, 46, 58]. Rather than using discrete representations as in the above methods, a recent line of work explores neural representations to encode the scene geometry and appearance as a continuous volumetric field [2, 24, 28, 31, 51]. We refer the readers to [54] for more comprehensive discussions of these works.

Many of these methods learn a mapping from 3D points and 2D viewing directions to color/density, which can be rendered using numerical integration – thus, rendering the color of a pixel for novel views requires hundreds of MLP evaluations. Instead, our approach directly learns the mapping from rays to color, and therefore supports efficient rendering with fewer network evaluations. In some ways, our method is similar to AutoInt [23], which predicts integrated radiance along ray segments and uses subdivision. The accuracy of AutoInt, however, is lower than NeRF due to the approximation of the nested volume rendering integral. By contrast, our approach leads to comparable performance to NeRF.

**Coordinate-based Representations.** Coordinate-based representations have emerged as a powerful tool for overcoming the limitations of traditional discrete representations (e.g., images, meshes, voxelized volumes). The core idea is to train an MLP to map an input coordinate to the desired target value such as pixel color [27, 34, 48, 65], signed distance [5, 35], occupancy [29], volume density [31], or semantic labels [70]. Like existing coordinate-based representation approaches, our method also learns the mapping from an input coordinate (ray) to a target scene property (color). Our core contributions lie in designing a coordinate embedding network that enables plausible view synthesis.

**Learned Coordinate Embedding.** Learned coordinate transformation or embedding has been applied to extend the capability of coordinate-based representations, such as NeRF, for handling higher dimensional interpolation problems. For example, several works tackle dynamic scene view synthesis [11, 22, 36, 56, 64] by mapping space at each time-step into a canonical frame. Others leverage coordinate embedding for modeling articulated objects such as the human body [25, 37]. Our ray-space embedding network can be viewed as locally deforming the input ray coordinates such that the MLP can produce plausible view interpolation while preserving faithful reconstruction of the source views.

**Subdivision.** Recent coordinate-based representations employ 3D space partitioning/subdivision either for improving the rendering efficiency [16, 41, 68] or representation accuracy [27, 40]. Our work also leverages spatial subdivision and shows improved rendering quality, particularly on scenes with large camera baselines.

## 3. Neural Light Fields

A neural radiance field [31] represents the appearance and geometry of a scene with an MLP $F_\theta : (\mathbf{x}_t, \vec{\omega}) \to (L_e(\mathbf{x}_t, \vec{\omega}), \sigma(\mathbf{x}_t))$ with trainable weights $\theta$. It takes as input a 3D position $\mathbf{x}_t$ and a viewing direction $\vec{\omega}$, and produces both the density $\sigma(\mathbf{x}_t)$ at point $\mathbf{x}_t$ and the radiance $L_e(\mathbf{x}_t, \vec{\omega})$ emitted at point $\mathbf{x}_t$ in direction $\vec{\omega}$.

One can generate views of the scene from this MLP using volume rendering:

$$L(\mathbf{x}, \vec{\omega}) = \int_{t_n}^{t_f} \underbrace{T(\mathbf{x}, \mathbf{x}_t)}_{\text{Acc. transmittance}} \underbrace{\sigma(\mathbf{x}_t)}_{\text{Density}} \underbrace{L_e(\mathbf{x}_t, \vec{\omega})}_{\text{Emitted radiance}} dt, \quad (1)$$

where $T(\mathbf{x}, \mathbf{x}_t) = e^{-\int_{t_n}^{t} \sigma(\mathbf{x}_s) ds}$ describes the accumulated transmittance for light propagating from position $\mathbf{x}$ to $\mathbf{x}_t$, for near and far bounds $t \in [t_n, t_f]$ of the scene. In practice, the integral is approximated using numerical quadrature:

$$L(\mathbf{x}, \vec{\omega}) \approx \sum_{k=1}^{N} \hat{T}(\mathbf{x}, \mathbf{x}_k)(1 - e^{-\sigma(\mathbf{x}_k)\Delta\mathbf{x}_k}) L_e(\mathbf{x}_k, \vec{\omega}), \quad (2)$$

where $\hat{T}(\mathbf{x}, \mathbf{x}_k) = e^{-\sum_{j=1}^{k} \sigma(\mathbf{x}_j)\Delta\mathbf{x}_j}$ and $\Delta\mathbf{x}_j = \mathbf{x}_j - \mathbf{x}_{j-1}$ is the distance between adjacent samples. An accurate approximation of (1) requires many samples and thus many neural network evaluations (on the order of hundreds) per ray.

**Light Fields.** Whereas a radiance field represents the radiance emitted at each point in space, a light field represents the total integrated radiance traveling along a ray. In other words, Equations (1) and (2) describe the relationship between the light field $L$ on the left-hand side and the radiance field $(L_e, \sigma)$ on the right-hand side. By optimizing a neural light field representation for $L$, given an input ray, one can predict this integrated radiance (and therefore render the color along this ray) with only a single neural network evaluation.
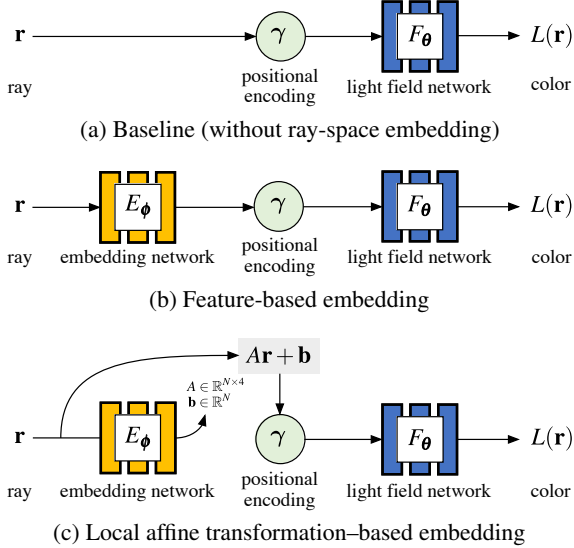
(a) Baseline (without ray-space embedding)



(b) Feature-based embedding



(c) Local affine transformation–based embedding

Figure 2. **Ray-space embedding networks**. Building upon the baseline (a), we explore two different schemes for ray-space embedding: (b) predicting the latent features and (c) predicting affine transformation parameters.

Assuming that we are operating on forward facing scenes, and that radiance remains constant along rays, we can construct a 4D ray space by intersecting rays with two planes, $\pi^{xy}$ and $\pi^{uv}$. We can then take the local plane coordinates, $(x, y) \in \pi^{xy}$ and $(u, v) \in \pi^{uv}$ of these ray intersections as the ray parameterization for the 4D light field, $L : \mathbf{r} = (x, y, u, v) \to \mathbf{c} = (r, g, b)$.

**Baseline Neural Light Fields.** To begin with, similar to how NeRF uses an MLP to represent radiance fields, we define an MLP $F_{\boldsymbol{\theta}}$ to represent a light field. It takes as input positionally encoded 4D ray coordinates $\boldsymbol{\gamma}(\mathbf{r})$ and outputs color $\mathbf{c}$ (integrated radiance) along each ray (Figure 2a):

$$L_{\text{base}}(\mathbf{r}) = F_{\boldsymbol{\theta}}(\boldsymbol{\gamma}(\mathbf{r})). \tag{3}$$

Unfortunately, this baseline approach is an unsatisfactory light field representation due to the following challenges. First, the captured input images only provide partial observations in the form of a sparse set of 2D slices of the full 4D ray space, so that each 4D ray coordinate in the input training data is observed at most once, assuming that no cameras share the same origin.

Second, light fields do not explicitly represent 3D scene geometry; hence the network a priori does not know how to interpolate the colors of unobserved rays from training observations. In other words, when querying a neural network representing a light field with unseen ray coordinates, multi-view consistency is not guaranteed. To address these challenges, we present two key techniques — ray-space embedding and subdivision — which substantially improve the rendering quality of the proposed neural light field representation.
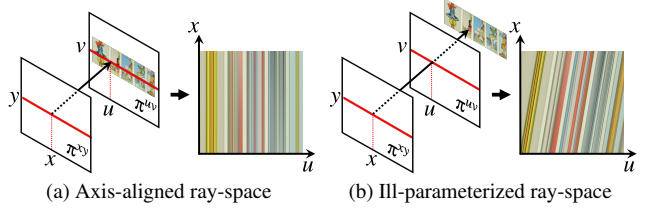


(a) Axis-aligned ray-space     (b) Ill-parameterized ray-space

Figure 3. **Importance of light field parameterization**. (a) It is easy to learn a light field for a well-parameterized ray space, but (b) a ray space with non-axis-aligned level sets makes the problem more challenging [53]. See Section 4 and appendix for discussions.

## 4. Ray-Space Embedding Networks

We first introduce *ray-space embedding networks*. These networks re-map the input ray-space into an embedded latent space. Intuitively, our goals for ray-space embedding networks are:

1. **Memorization:** Map disparate coordinates in input 4D ray-space that observe the same 3D point to the same location in the latent space. This allows for more stable training and better allocation of network capacity, and thus better memorization of input views.

2. **Interpolation:** Produce an interpolable latent space, such that querying unobserved rays preserves multi-view consistency and improves view synthesis quality.

**Feature-Space Embedding.** A straightforward approach would be to learn a nonlinear mapping from 4D ray coordinates to a latent feature space via an MLP, $E_{\boldsymbol{\phi}}^{\text{feat}} : \mathbf{r} \to \mathbf{f} \in \mathbb{R}^N$, where $N$ is the feature space dimension (see Figure 2b). This embedding network can produce a nonlinear, many-to-one mapping from distinct ray coordinates into shared latent features, which would allow for better allocation of capacity in the downstream light field network $F_{\boldsymbol{\theta}}$. The finite capacity of $F_{\boldsymbol{\theta}}$ incentivizes this "compression" and thus encourages rays with similar colors to be mapped to nearby features in the latent space; effectively a form of implicit triangulation.

After embedding, we then feed the positionally encoded $N$-dimensional feature into the light field MLP $F_{\boldsymbol{\theta}}$:

$$L_{\text{feat}}(\mathbf{r}) = F_{\boldsymbol{\theta}}(\boldsymbol{\gamma}(E_{\boldsymbol{\phi}}^{\text{feat}}(\mathbf{r}))). \tag{4}$$

While this approach compresses and helps to find correspondences in ray space, it is not enough to facilitate the interpolation between these correspondences (see Section 6).

**Light Field Re-parameterization.** To understand our goals better, consider a scene containing an axis-aligned textured square at a fixed $z$-depth. Again, we assume a two-plane parameterization for the input ray-space of this scene. If one of the planes is at the depth of the square, then the light field will only vary in two coordinates, as illustrated in Figure 3a. In other words, the light field is only a function of $(u, v)$ and is constant in $(x, y)$. Since positional encoding is applied separately for each input coordinate, the baseline

model (3) can interpolate the light field perfectly if it chooses to ignore the positionally encoded coordinates $(x, y)$.

On the other hand, if the two planes in the parameterization do not align with the textured square, the color level sets (e.g., the line structures in the $ux$ slices) of the light field are 2D affine subspaces that are not axis-aligned (see Figure 3b). Axis-aligned positional encoding, as used by NeRF [31], will yield interpolation kernels [53] that are ill-suited for capturing these sheared subspaces, especially from incomplete observations. One can therefore frame the task of effective interpolation as learning an optimal re-parameterization of the light field for a scene, such that the 2D color level sets for each 3D scene point are axis-aligned. See the appendix for a more extensive discussion.

**Local Affine-Transformation Embedding.** Given this intuition, we describe our architecture for ray-space embedding, which learns local affine transformations for each coordinate in ray space. We learn local affine transformations rather than a single global transformation because different rays may correspond to points with different depths, which will cause the shape of color level sets to vary.

To this end, we employ an MLP $E_\phi : \mathbf{r} \to (A \in \mathbb{R}^{N \times 4}, \mathbf{b} \in \mathbb{R}^N)$. The output of the network is a $N \times 4$ matrix $A$, as well as an $N$-vector $\mathbf{b}$ representing a bias (translation), which together form a 4D$\to N$D affine transformation. This affine transformation is applied to the input ray coordinate $\mathbf{r}$ before being positionally encoded and passed into the light field network $F_\theta$:

$$L(\mathbf{r}) = F_\theta(\gamma(A\mathbf{r} + \mathbf{b})), \text{ where } (A, \mathbf{b}) = E_\phi(\mathbf{r}). \quad (5)$$

See Figure 2c for an illustration of this model.

While setting $N = 4$ above is perfectly reasonable, we use $N = 32$ in practice, which can be seen as providing multiple possible re-parameterizations per ray. Also note that learning arbitrary affine transforms, rather than a single $z$-depth for each ray, allows the network to better capture both *angular frequencies* in the light field (due to object depth), as well as *spatial frequencies* (due to object texture).

## 5. Subdivided Neural Light Fields

Although this approach works well for dense light fields, our embedding network struggles to find long-range correspondences between training rays when training data is too sparse. Moreover, even if it can discover correspondences, interpolation for unobserved rays in between these correspondences remains underconstrained.

To resolve these issues, we propose learning a voxel grid of local light fields that covers the entire 3D scene. This approach is motivated by the following observation: if we parameterize a local light field for a voxel by intersecting rays with the voxel's front and back planes (See Figure 4), the color level sets of the local light field are already almost
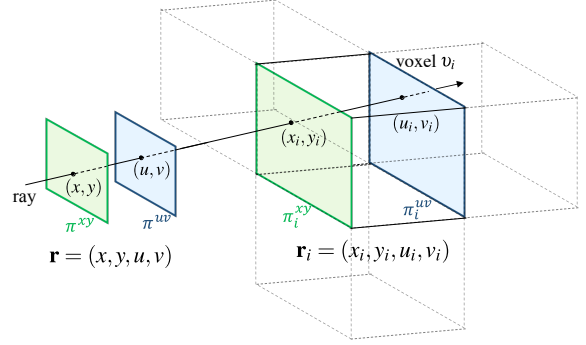


Figure 4. **Re-parameterization for local light fields**. A ray $\mathbf{r} = (x, y, u, v)$ defined globally is re-parameterized with respect to each local light field in voxel $v_i$, by intersecting the front and back faces of $v_i$. This yields local ray coordinates $\mathbf{r}_i = (x_i, y_i, u_i, v_i)$.

axis-aligned. As such, the ray-space embedding for each voxel can be simple, only inducing small shears of the input ray-space to allow for good interpolation. On the other hand, a global light field requires learning a complex ray-space embedding that captures occlusions and significant depth changes.

Therefore a local light field can be learned much more easily than an entire global light field with the same training data. This still requires that we solve an assignment problem: we must know which rays to use to train each local light field. While we can exclude all rays that do not intersect the voxel containing the local light field, many other rays may be occluded before they hit the voxel. Therefore, these rays should also be excluded during training.

A simple way to handle ray assignment is by learning (per-ray) opacity. If the opacity accumulated along a ray before it hits a voxel is high, the ray should receive a little contribution from this voxel. We, therefore, modify our light field network also to produce integrated opacity, or alpha for each ray. It is important to note that the opacity can vary depending on which ray is queried through a particular voxel — opacity is a function of the ray, not the voxel itself.

**Subdivided Volume Rendering.** We detail how a set of local light fields can be rendered to form new images. Given a voxel grid in 3D space, we place a local light field within each voxel. A voxel's light field is parameterized with respect to its front and back planes (See Figure 4), yielding ray coordinates $\mathbf{r}_i$, where $i$ is some voxel index, for any ray $\mathbf{r}$ defined globally. If 3D space is subdivided into $M$ voxels, there exist $M$ distinct parameterizations $\{\mathbf{r}_i\}_{i=1}^M$ of each $\mathbf{r}$.

To implement this model, the embedding network $E_\phi$ is augmented such that it takes the positionally encoded voxel index $\gamma(i)$ in addition to the 4D ray coordinates, i.e., $E_\phi : (\mathbf{r}_i, \gamma(i)) \to (A_i, \mathbf{b}_i)$. The light field network $F_\theta$ is augmented in the same way and predicts both color *and* alpha: $F_\theta : (\gamma(A_i\mathbf{r}_i + \mathbf{b}_i), \gamma(i)) \to (\mathbf{c}_i, \alpha_i)$. Again, we emphasize that $\alpha_i$ is a function of the embedded ray, not just the voxel.
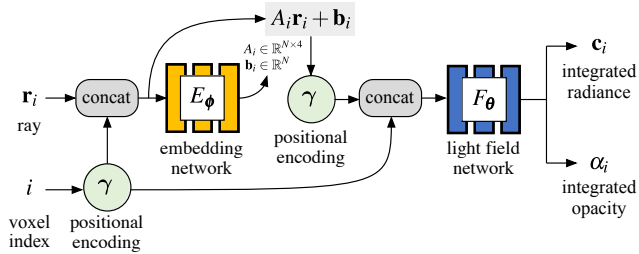
Figure 5. **Local light fields via spatial subdivision**. Each local light field takes a locally re-parameterized ray $\mathbf{r}_i$ and its voxel index $i$, and predicts the radiance and opacity *both* integrated for the ray segment traveling through the voxel. Both the embedding network $E_\phi$ and the light field network $F_\theta$ are conditioned on positionally encoded voxel index $\gamma(i)$.

Given this model, rendering works in the following way. First, the set of voxels $\mathcal{V}(\mathbf{r})$ that the ray $\mathbf{r}$ intersects are identified. The ray is then intersected with each voxel's front and back planes $\pi_i^{xy}$ and $\pi_i^{uv}$. This process yields a series of local ray coordinates $\{\mathbf{r}_i\}_{v_i \in \mathcal{V}(\mathbf{r})}$. The color $\mathbf{c}_i$ and alpha $\alpha_i$ of each $\mathbf{r}_i$ are then computed as:

$$(\mathbf{c}_i, \alpha_i) = F_\theta(\gamma(A_i\mathbf{r}_i + \mathbf{b}_i), \gamma(i)), \qquad (6)$$

where the local affine-transformation is obtained from the embedding network:

$$(A_i, \mathbf{b}_i) = E_\phi(\mathbf{r}_i, \gamma(i)). \qquad (7)$$

The final color of ray $\mathbf{r}$ is then over-composited:

$$\mathbf{c} = \sum_{i \in \mathcal{V}(\mathbf{r})} \Big( \prod_{j \in \mathcal{V}(\mathbf{r}) \wedge j < i} (1 - \alpha_j) \Big) \alpha_i \mathbf{c}_i, \qquad (8)$$

which assumes that we sort the voxels by their distance from the ray origin in ascending order (see Figure 5).

# 6. Experiments

**Datasets.** We evaluate our model in two angular sampling regimes: sparse (wide-baseline) and dense (narrow-baseline), using the Real Forward-Facing dataset [31], the Shiny dataset [61], and the Stanford Light Field dataset [60]. Each scene from the Real Forward-Facing and Shiny datasets consist of a few dozen images with relatively large angular separation in between views. We followed their authors' train/test split and hold out every 8th image from training.

The Stanford light fields are all captured as $17 \times 17$ dense 2D grids of images. We take every 4th image in *both* horizontal and vertical directions and use the resulting $5 \times 5$ subsampled light field for training (thus using 25 images out of 289), while reserving the rest of the images for testing.

We downsample the images in the Real Forward-Facing and Shiny datasets to $504 \times 378$ pixels, and use half the original resolution of the Stanford images. For quantitative metrics, we report PSNR, SSIM, LPIPS, and FPS; where FPS is computed for $512 \times 512$ pixel images.

Table 1. **Quantitative comparisons.** We show the aggregated metrics over all scenes in each of the following datasets: Real Forward-Facing [31], *Undistorted* Real Forward-Facing [61], Shiny [61], and Stanford [60].

| Dataset | Method | PSNR↑ | SSIM↑ | LPIPS↓ | FPS↑ |
|---|---|---|---|---|---|
| RFF | NeRF [31] | **27.928** | **0.916** | 0.065 | 0.07 |
| | AutoInt [23] | 25.531 | 0.853 | 0.156 | 0.26 |
| | Ours (subdivision) | 27.454 | 0.905 | **0.060** | **0.34** |
| Undistorted RFF | NeX [61] | **27.992** | **0.924** | **0.052** | 0.13 |
| | Ours (subdivision) | 27.355 | 0.905 | 0.059 | **0.34** |
| Shiny | NeX [61] | 28.268 | **0.940** | **0.036** | 0.13 |
| | Ours (subdivision) | **28.499** | 0.931 | 0.038 | **0.34** |
| Stanford | NeRF [31] | 37.559 | 0.979 | 0.037 | 0.07 |
| | NeX [61] | 36.286 | 0.972 | 0.033 | 0.13 |
| | X-Fields [2] | 35.559 | 0.976 | 0.036 | **20.00** |
| | Ours (*no* subdivision) | **38.054** | **0.982** | **0.020** | 9.14 |

**Baseline Methods.** We compare our model to NeRF [31], AutoInt [23] (with 32 sections), and NeX [61] on the Real Forward-Facing dataset [31]; and additionally to NeX on the Shiny dataset [61]. As in [61], we perform evaluation on NeX before baking. More information about this is provided in the appendix. In Table 1, we include metrics for both the Real Forward-Facing and *Undistorted* Real Forward-Facing dataset, a variant of the original dataset that is supported by the NeX codebase. We did not include the following methods, although they are related, for various reasons: PlenOctrees [68], KiloNeRF [41], and NSVF [24] operate on bounded 360-degree scenes only; SNeRG [16] strictly underperforms NeRF on real forward-facing scenes.

On Stanford [60], we compare our model to NeRF, NeX, and X-Fields [2], which is a state of the art method in image-based light field interpolation. Kalantari et al. [17] would also be a relevant comparison, but X-Fields outperforms it, hence we did not include it in our evaluations.

**Subdivision.** For small-baseline light fields (Stanford [60]), *we do not use subdivision*, and use a single-evaluation per ray model with ray-space embedding. For sparse light fields (Real Forward Facing [31] and Shiny [61]) we place local light fields within each voxel in a $32^3$ regular voxel grid that covers all of NDC space, with the exception of the CD and Lab sequences in the Shiny dataset, where we use a coarser $4^3$ voxel grid.

**Implementation Details.** For our *feature-based embedding* (4), we set the embedding space dimension to $N = 32$. The embedded feature $\mathbf{f}$ is $\ell_2$-normalized and then multiplied by $\sqrt{32}$, such that on average, each output feature has a square value close to one. For our *local affine transformation–based embedding* (5), we infer $32 \times 4$ matrices and 32D bias vectors. We then normalize the matrices $A$ with their Frobenius norm and multiply them by $\sqrt{32 \cdot 4}$. We use a *tanh*-activation for our predicted bias vectors.

For positional encoding $\gamma(\cdot)$ we use $L = 8$ frequency bands for our subdivided model and $L = 10$ frequency bands for our one-evaluation model. Higher frequencies are gradu-
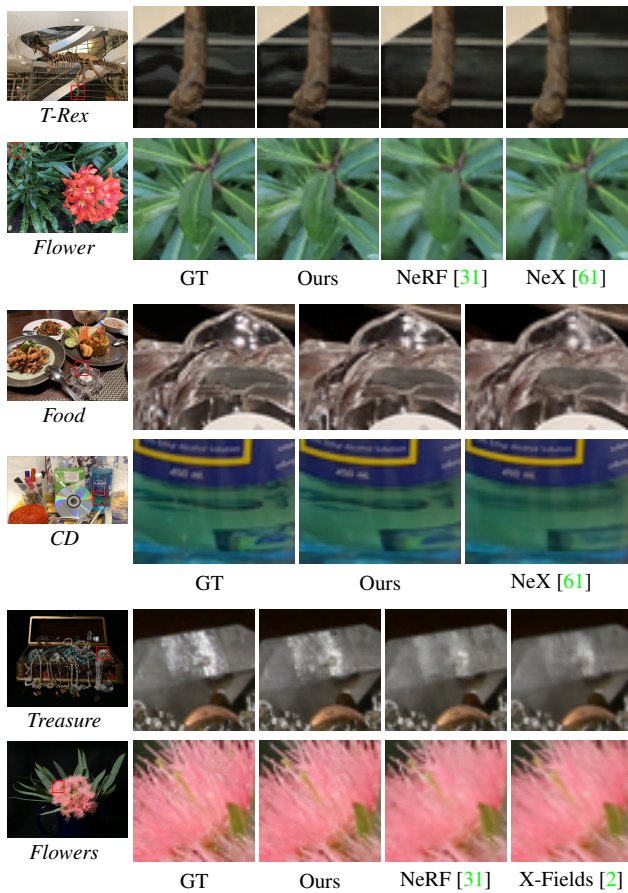
Figure 6. **Visual comparisons.** Close-ups reveal our representation leads to more faithful reconstruction of view-dependent effects, particularly distorted reflections and refractions (*Horns* and *T-Rex* from RFF; *Treasure* from Stanford; and *Food* and *CD* from Shiny) as well as sharper details (*Flower* from RFF and *Flowers* from Stanford). See our project website for more results.

ally eased in using windowed positional encoding, following Park et al. [36] — over 50k iterations with subdivision, or 80k without. For both the embedding and color models we use a 8-layer, 256 hidden unit MLP with one skip connection. We share a single network each for all local light fields for both the embedding $E_\phi$ and color $F_\theta$ in our subdivided representation, which are "indexed" (see section 5) with the voxel center points $\mathbf{p}_i$, i.e., $i \equiv \mathbf{p}_i$.

We use a batch size of 4,096 for our subdivided model, and a batch size of 8,192 for our one-evaluation model during training and run all experiments on a single Tesla V100 GPU with 16 GB of memory. Please refer to the appendix and our project webpage for further details and results. We will release the source code and pre-trained models to facilitate future research.

## 7. Results

**Quantitative Evaluation.** Table 1 compares view synthesis quality using the metrics aggregated over held-out im-

ages in all scenes. For the Real Forward-Facing dataset, our method outperforms AutoInt in both quality ($\sim$2 dB) and speed ($\sim$30%) while performing competitively with NeRF on PSNR and SSIM (within 0.5 dB in quality) and rendering more than $\times$4.5 times faster. We also outperform NeRF on LPIPS. Additionally, we are competitive with NeX on the *Undistorted* Real Forward-Facing dataset (within 0.65 dB in quality). For the Shiny dataset, we outperform NeX in terms of PSNR, while performing slightly worse on SSIM and LPIPS metrics. We do especially well on the *CD* and *Lab* scenes with challenging view dependence, outperforming NeX by > 3dB each.

On Stanford, our method outperforms NeRF, NeX, and X-Fields in all three metrics. Specifically, our model performs better on more complex scenes, e.g., *Treasure* and *Bulldozer*, with complicated geometry and high-frequency textures, while NeRF performs slightly better on simpler scenes like *Chess* and *Beans*. Our model renders a frame in about 0.11 seconds in this setup, which is faster than all other methods apart from X-Fields.

**Training Time.** For both RFF and Shiny, it takes about 20 hours to train our model at $32^3$ resolution, while for the Stanford data, we train our model for about 10 hours. NeRF takes approximately 18 hours to train on all datasets, while NeX takes approximately 36 hours on 2 GPUs. For X-Fields, we double the default capacity of their model as well as the number of images used for interpolation; and train X-Fields for 8 hours rather than the maximum time of 172 minutes listed in their paper [2]. See the appendix for more details.

**Reconstruction of View Dependence.** Our approach reconstructs view dependence better than the baseline methods in both sparse and dense regimes (see Figures 1(*left*) and 6), due to the flexibility of the embedding network to predict different affine transforms for different regions of ray-space. This effectively produces locally deforming depth estimates, which can better model reflected and refracted content, for example. On the other hand, unlike our method, volumetric representations use a fixed global geometry, and are either late-conditioned on viewing direction [31], or use lower frequency basis functions for view dependence [61]. Thus, they struggle to represent high-frequency view dependent appearance that breaks multi-view color constraints.

In general, reproducing complex view dependence is a major challenge for volumetric representations, and is one of the primary contributions of our work. The difference is especially prominent in the Shiny dataset comparisons in Figure 1, as well as the results for *T-rex*, *Food*, *CD*, and *Treasure* in Figure 6.

**Memory Footprint.** Table 2 shows the model sizes. Our model has a similar memory footprint to NeRF. NeX explicitly represents albedo for 16 MPI layers, leading to larger memory consumption than pure coordinate-based neural representations. Overall, our method provides the best trade-off

Table 2. **Memory footprint** of our model compared to other methods in terms of the number of trainable weights in each network. For X-Fields [2] we include the memory footprint of the model, and of the training images which must be maintained in order to synthesize novel views.

| Method | Ours | NeRF [31] | NeX [61] | AutoInt [23] | X-Fields [2] |
|--------|------|-----------|----------|--------------|--------------|
| Size (MB) | 5.4 | 4.6 | 38.4 | 4.6 | 3.8 / 75 |

Table 3. **Ablations on the ray-embedding network** with subdivision (on Real Forward-Facing [31]) and without subdivision (on Stanford [60] light fields).

| Subdivision | Method | Embedding | PSNR↑ | SSIM↑ | LPIPS↓ |
|-------------|--------|-----------|-------|-------|--------|
| Used (On RFF) | Ours | Not used (3) | 26.696 | 0.888 | 0.079 |
| | | Feature (4) | 26.922 | 0.891 | 0.073 |
| | | Local affine (5) | 27.454 | 0.905 | **0.060** |
| | NeRF [31] | | **27.928** | **0.916** | 0.065 |
| Not used (On Stanford) | Ours | Not used (3) | 25.111 | 0.841 | 0.098 |
| | | Feature (4) | 37.120 | 0.9792 | 0.030 |
| | | Local affine (5) | **38.054** | **0.982** | **0.020** |
| | NeRF [31] | | 37.559 | 0.9790 | 0.037 |

Table 4. **Ablations on the subdivision** on sparse light fields (Real Forward-Facing [31]).

| Method | Subdivision | PSNR↑ | SSIM↑ | LPIPS↓ | FPS↑ |
|--------|-------------|-------|-------|--------|------|
| Ours | $4^3$ grid | 25.579 | 0.867 | 0.094 | **2.82** |
| | $8^3$ grid | 26.860 | 0.893 | 0.070 | 1.42 |
| | $16^3$ grid | 27.350 | 0.904 | 0.062 | 0.69 |
| | $32^3$ grid | 27.454 | 0.905 | **0.060** | 0.34 |
| AutoInt [23] | 8 sections | 24.136 | 0.820 | 0.176 | 0.94 |
| | 16 sections | 24.898 | 0.836 | 0.167 | 0.51 |
| | 32 sections | 25.531 | 0.853 | 0.156 | 0.26 |
| NeRF [31] | | **27.928** | **0.916** | 0.065 | 0.07 |

between quality, rendering speed, and memory among the baselines we compared to; see the graphs in Figure 1(*right*).

**Ablations on Embedding Networks.** We study the effect of our proposed embedding networks via ablations. See Table 3 for quantitative evaluation with varying configurations.

With no embedding and only positional encoding (baseline), the network does not produce multi-view consistent view synthesis. The use of our feature-based embedding leads to a boost in quality over the baseline. In particular, it achieves our first objective of "memorization" (identifying rays that observe the same 3D point), but still struggles with view interpolation. Our local affine transformation embedding achieves better interpolation quality because it implicitly encourages color level sets to be locally linear. Thus, we learn a model that interpolates in a (locally) view consistent way for free.

**Ablations on Subdivision.** Table 4 and Figure 1(*right*) summarize the trade-off between the quality and rendering speed of our model on various resolutions of spatial subdivision, compared to NeRF and AutoInt. For each subdivision level, we use a regular $N^3$ voxel grid, where $N \in (4, 8, 16, 32)$. As expected, our model archives better quality view synthesis as we employ more local light fields with finer spatial subdivision, trading off the rendering speed. The render-time increases linearly in $N$. Our method provides better quality and speed than AutoInt at all comparable subdivision levels.

**Limitations.** One limitation of our work is that we focus on light fields parameterized with two planes. Although common, this does not encompass 360° scenes. In the appendix, we show preliminary results for our method on 360° scenes using Plücker coordinates [49], which again outperforms Au-

toInt. However, color level sets for Plücker coordinates are no longer affine. Thus, we believe that improving the design of embedding networks for different parameterizations could lead to a larger boost in performance.

Additionally, our representation requires subdivision for sparse light fields. This comes at the cost of both increased render time and training time. We hope to investigate other network designs or regularization schemes, which encourage the representation to be multi-view consistent without subdivision. Adaptive subdivision [24, 27] is another interesting direction for future work which could lead to better quality without sacrificing speed.

# 8. Conclusion

We present a novel ray-space embedding approach for learning neural light fields, which achieves state-of-the-art quality on small-baseline datasets. To better handle sparse input, we leverage spatial subdivision with a voxel grid of local light fields, which improves quality at the cost of increased render time. Our subdivided representation enables comparable performance to existing models, and achieves a better trade-off between quality, speed, and memory. Additionally, in both regimes, our method can handle complex view dependent effects that existing state-of-the-art volumetric scene representations do not faithfully reproduce [31, 61].

We believe that our approach is orthogonal to contemporary works that leverage hybrid or explicit grid-based representations for view synthesis [32, 52, 67], as well as approaches that use many small MLPs within a 3D volume [41]. As such, our approach can potentially be used to reduce the required grid resolution and number of sample points per ray for these methods — which may further improve rendering quality, speed, and memory overhead . We hope that our work can open up new avenues for view synthesis; and further, that its insights can be leveraged for other signal-representation problems (such as image/video representation) or high-dimensional interpolation problems (such as view-time or view-time-illumination interpolation [2]).

# References

[1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *ECCV*, 2020. 3

[2] Mojtaba Bemana, Karol Myszkowski, Hans-Peter Seidel, and Tobias Ritschel. X-fields: implicit neural view-, light- and time-image interpolation. *ACM Trans. Graph.*, 39(6):257:1–257:15, 2020. 3, 6, 7, 8, 14

[3] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew Duvall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh representation. *ACM Transactions on Graphics (TOG)*, 39(4):86–1, 2020. 3

[4] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 2001. 2, 3

[5] Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *ECCV*, 2020. 3

[6] Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM Transactions on Graphics (TOG)*, 32(3):1–12, 2013. 3

[7] Abe Davis, Marc Levoy, and Fredo Durand. Unstructured light fields. In *Eurographics*, 2012. 2

[8] William Falcon and The PyTorch Lightning team. PyTorch Lightning, 3 2019. 11

[9] Brandon Yushan Feng and Amitabh Varshney. Signet: Efficient neural representation for light fields. In *ICCV*, 2021. 2, 3

[10] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world's imagery. In *CVPR*, 2016. 3

[11] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, 2021. 3

[12] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. 2

[13] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. 2

[14] Peter Hedman, Suhib Alsisan, Richard Szeliski, and Johannes Kopf. Casual 3d photography. *ACM Transactions on Graphics (TOG)*, 36(6):1–15, 2017. 3

[15] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (TOG)*, 37(6):1–15, 2018. 3

[16] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. *ICCV*, 2021. 2, 3, 6

[17] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Transactions on Graphics (TOG)*, 35(6):1–10, 2016. 3, 6

[18] Petr Kellnhofer, Lars Jebe, Andrew Jones, Ryan Spicer, Kari Pulli, and Gordon Wetzstein. Neural lumigraph rendering. In *CVPR*, 2021. 2

[19] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. *Computer Graphics Forum*, 40(4):29–43, 2021. 3

[20] Johannes Kopf, Kevin Matzen, Suhib Alsisan, Ocean Quigley, Francis Ge, Yangming Chong, Josh Patterson, Jan-Michael Frahm, Shu Wu, Matthew Yu, et al. One shot 3d photography. *ACM Transactions on Graphics (TOG)*, 39(4):76–1, 2020. 3

[21] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. 2

[22] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 3

[23] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *CVPR*, 2021. 2, 3, 6, 8, 11, 13, 14

[24] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 3, 6, 8

[25] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Transactions on Graphics (TOG)*, 2021. 3

[26] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. In *SIGGRAPH*, 2019. 3

[27] Julien N. P. Martel, David B. Lindell, Connor Z. Lin, Eric R. Chan, Marco Monteiro, and Gordon Wetzstein. Acorn: Adaptive coordinate networks for neural scene representation. *ACM Transactions on Graphics (TOG)*, 40(4), 2021. 3, 8

[28] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 3

[29] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, 2019. 3

[30] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 3

[31] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 3, 5, 6, 7, 8, 11, 13, 14

9

[32] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 2, 8

[33] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of compact neural radiance fields using depth oracle networks. *Comput. Graph. Forum*, 40(4):45–59, 2021. 2

[34] Michael Oechsle, Michael Niemeyer, Christian Reiser, Lars Mescheder, Thilo Strauss, and Andreas Geiger. Learning implicit surface light fields. In *International Conference on 3D Vision (3DV)*, 2020. 3

[35] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 3

[36] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021. 3, 7

[37] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 3

[38] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)*, 36(6):1–11, 2017. 3

[39] Chen Quei-An. Nerf_pl: a pytorch-lightning implementation of nerf, 2020. 11

[40] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. In *CVPR*, 2021. 3

[41] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, 2021. 2, 3, 6, 8

[42] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, 2020. 3

[43] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *CVPR*, 2021. 3

[44] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *arXiv preprint arXiv:2110.06635*, 2021. 3

[45] Lixin Shi, Haitham Hassanieh, Abe Davis, Dina Katabi, and Fredo Durand. Light field reconstruction using sparsity in the continuous fourier domain. *ACM Transactions on Graphics (TOG)*, 34(1):1–13, 2014. 3

[46] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *CVPR*, 2020. 3

[47] Harry Shum and Sing Bing Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing 2000*, volume 4067, pages 2–13, 2000. 3

[48] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *NeurIPS*, 2020. 3

[49] Vincent Sitzmann, Semon Rezchikov, William T Freeman, Joshua B Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. In *NeurIPS*, 2021. 2, 3, 8

[50] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 3

[51] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 3

[52] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. *arXiv preprint arXiv:2111.11215*, 2021. 2, 8

[53] Matthew Tancik, Pratul P Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *NeurIPS*, 2020. 4, 5, 12

[54] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. State of the art on neural rendering. *Computer Graphics Forum*, 39(2):701–727, 2020. 1, 3

[55] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, Tomas Simon, Christian Theobalt, Matthias Niessner, Jonathan T. Barron, Gordon Wetzstein, Michael Zollhoefer, and Vladislav Golyanik. Advances in neural rendering, 2021. 1

[56] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhofer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 3

[57] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *CVPR*, 2020. 3

[58] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3d scene inference via view synthesis. In *ECCV*, 2018. 3

[59] Suren Vagharshakyan, Robert Bregovic, and Atanas Gotchev. Light field reconstruction using shearlet transform. *IEEE TPAMI*, 40(1):133–147, 2017. 3

[60] Bennett Wilburn, Neel Joshi, Vaibhav Vaish, Eino-Ville Talvala, Emilio R. Antúnez, Adam Barth, Andrew Adams, Mark Horowitz, and Marc Levoy. High performance imaging using large camera arrays. *ACM Trans. Graph.*, 24(3):765–776, 2005. 2, 6, 8, 11, 12, 14

[61] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. NeX: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2, 3, 6, 7, 8, 11, 14

[62] Gaochang Wu, Yebin Liu, Qionghai Dai, and Tianyou Chai. Learning sheared epi structure for light field reconstruction. *IEEE TIP*, 28(7):3261–3273, 2019. 3

[63] Gaochang Wu, Mandan Zhao, Liangyong Wang, Qionghai Dai, Tianyou Chai, and Yebin Liu. Light field reconstruction using deep convolutional network on epi. In *CVPR*, 2017. 3

[64] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *CVPR*, 2021. 3

[65] Fanbo Xiang, Zexiang Xu, Milos Hasan, Yannick Hold-Geoffroy, Kalyan Sunkavalli, and Hao Su. Neutex: Neural texture mapping for volumetric neural rendering. In *CVPR*, 2021. 3

[66] Henry Wing Fung Yeung, Junhui Hou, Jie Chen, Yuk Ying Chung, and Xiaoming Chen. Fast light field reconstruction with deep coarse-to-fine modeling of spatial-angular clues. In *ECCV*, 2018. 3

[67] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021. 2, 8

[68] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *CVPR*, 2021. 2, 3, 6

[69] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *CoRR*, abs/2010.07492, 2020. 2

[70] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew Davison. In-place scene labelling and understanding with implicit scene representation. In *ICCV*, 2021. 3

[71] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM Transactions on Graphics (TOG)*, 2018. 3

## A. Implementation Details

Our code is written entirely in python, using the PyTorch Lightning framework [8]. The design of our codebase was inspired by [39]. Below we include additional implementation details for our method for both the dense dataset (Stanford Light Field [60]) and sparse forward facing datasets (Real Forward-Facing [31] and Shiny [61]). We also include information about the NeRF [31], NeX [61], and AutoInt [23] baselines.

### A.1. Stanford Light Field Dataset

Calibration information is not provided with the Stanford Light Field dataset, apart from the $(x, y)$ positions of all images on the camera plane $\pi^{xy}$. All Stanford Light Fields are parameterized with respect to a plane $\pi^{uv}$ that approximately cuts through the center of each scene. Thus, a pixel coordinate in an image corresponds the location that a ray intersects $\pi^{uv}$. We heuristically set the location of $\pi^{xy}$ to $z = -1$, and the location of the object plane $\pi^{uv}$ to $z = 0$. We scale the camera positions so that they lie between $[-0.25, 0.25]$ in both $x$ and $y$, and the pixel coordinates on $\pi^{uv}$ so they lie between $[-1, 1]$. The camera positions, and the vector from

the camera origin to the location of the pixel on the object plane then comprise our ray origins and directions.

For our method, we take camera coordinates and object plane coordinates, which correspond to intersections of the rays with $\pi^{xy}$ and $\pi^{uv}$, as our initial two-plane ray parameterization.

For NeRF/NeX, we set the near distance to 0.5 and the far distance to 2 for all scenes, except for the Knights scene, where we set the near distance to 0.25. Note that, as defined, the scene coordinates may differ from the true (metric) world space coordinates by a projective transform. However, multi-view constraints still hold (intersecting rays remain intersecting, epipolar lines remain epipolar lines), and thus NeRF is able to learn an accurate volume.

### A.2. Real Forward-Facing Dataset

For the Real Forward-Facing Dataset, we perform all experiments in NDC space. For our subdivided model, we reparameterize each ray **r** within a voxel $\upsilon$ by first transforming space so that the voxel center lies at the origin. We then intersect the transformed ray with the voxel's front and back planes, and take the *xy* coordinates of these intersections, which lie within the range $[-voxel\ width, +voxel\ width]$ in all dimensions as the parameterization.

### A.3. Shiny Dataset

Our procedure for evaluating on Shiny [61] is identical to the Real Forward-Facing dataset. We perform experiments in NDC space and use a $32^3$ voxel grid that covers all of NDC space on all scenes *except* CD/Lab where we use a coarser $4^3$ grid.

### A.4. NeRF and AutoInt Baselines

For the Stanford Dataset, we train NeRF for 400k iterations with a batch size of 1,024 for all scenes. For the Real Forward-Facing dataset, AutoInt does not provide pretrained models and training with their reported parameters on a V100 GPU with 16GB of memory leads to out-of-memory errors. They also do not provide multi-GPU training code. As such, we report the quantitative metrics for their method and for NeRF published in their paper, which come from models trained at the same resolution ($504 \times 378$), and using the same heldout views as ours on the Real Forward-Facing Dataset.

### A.5. NeX Baseline

We train NeX on all datasets (Stanford, Undistorted RFF, Shiny), using their public codebase, for 4000 epochs on all scenes (the default number of epochs in their training script), or about 36 hours each. We use their multi-GPU training code to split training over 2 16GB V100 GPUs. We note that NeX can perform real-time rendering after discretizing their view dependent basis functions into $400 \times 400$ textures.

However the NeX codebase does not include evaluation code for their real-time renderable MPIs, and thus the numbers for NeX in Table 1 of the main paper are reported *before* baking. While this presumably leads to an increase in quality, it takes a longer time to render images, hence the smaller FPS scores in Table 1.

## B. Importance of Light Field Parameterization

Here, we expand on the discussion in Section 4 of the main paper and describe why our light field parameterization is crucial for enabling good view interpolation. Let the two planes in the initial two plane parameterization be $\pi^{xy}$ and $\pi^{uv}$, with local coordinates $(x, y)$ and $(u, v)$. In addition, let us denote the plane of the textured square as $\pi^{st}$ with local coordinates $(s, t)$, and assume that it is between $\pi^{xy}$ and $\pi^{uv}$. Assume, without loss of generality, that the depth of $\pi^{xy}$ is 0, and suppose the depths of $\pi^{st}$ and $\pi^{uv}$ are $z_{st}$ and $z_{uv}$ respectively.

For a ray originating at $(\hat{x}, \hat{y})$ on $\pi^{xy}$ and passing through $(\hat{s}, \hat{t})$ on $\pi^{st}$, we can write (by similar triangles):

$$\hat{s} - \hat{x} = (\hat{u} - \hat{x})\frac{z_{st}}{z_{uv}}, \tag{9}$$

$$\hat{t} - \hat{y} = (\hat{v} - \hat{y})\frac{z_{st}}{z_{uv}}, \tag{10}$$

which gives

$$\hat{u} = \frac{\hat{x}(z_{st} - z_{uv}) + \hat{s}}{z_{st}}, \tag{11}$$

$$\hat{v} = \frac{\hat{y}(z_{st} - z_{uv}) + \hat{t}}{z_{st}}. \tag{12}$$

Recall that the positional encoding of the 4D input parameterization $\gamma(\hat{x}, \hat{y}, \hat{u}, \hat{v})$ will be fed into the light field network. Thus, the network will produce interpolation kernels aligned with $\hat{u}$ and $\hat{v}$. However for perfect interpolation, we would like the output of the light field network to only depend on $(\hat{s}, \hat{t})$, or for the interpolation kernels to be aligned with $(\hat{s}, \hat{t})$. It can be observed in equations (3) and (4) that the greater the distance $(z_{st} - z_{uv})$, the larger the difference between $(\hat{u}, \hat{v})$ and $(\hat{s}, \hat{t})$, and the *less aligned* the interpolation kernels become with $(\hat{s}, \hat{t})$.

On the other hand, by learning a re-parameterization of the light field, such that $\pi^{uv}$ is moved towards $\pi^{st}$ (i.e. reducing the distance $(z_{st} - z_{uv})$), we align the color network's interpolation kernels with $(\hat{s}, \hat{t})$. As in the feature-embedding approach, the finite capacity of the light field MLP will drive the embedding network to learn to map rays intersecting the same point on the textured square to the same point in the latent space—and thus will drive learning of an optimal re-parameterization that leads to good interpolation.

Table B.1. **Empirical validation of parameterization.** We calculate PSNR, SSIM, and LPIPS for our model without embedding trained on different initial parameterizations.

| $\pi^{uv}$ location | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| $z = 0$ | **29.631** | **0.937** | **0.059** |
| $z = 1$ | <u>24.411</u> | <u>0.848</u> | <u>0.096</u> |
| $z = 3$ | 21.491 | 0.790 | 0.146 |

### B.1. Empirical Validation

We claim that quality of a *baseline* neural light field is correlated with $(z_{st} - z_{uv})$. In order to support this claim, we perform a set of simple experiments with parameterization. In particular, we choose the *Amethyst* scene from the Stanford Light Field dataset [60], which has very little depth variation. As described in Section A.1, each Stanford light field has the object plane $\pi^{st}$ at $z = 0$. We re-parameterize the input light field for $\pi^{uv}$ at $z = 0, 1, 3$, train our model *without* an embedding network, and report validation metrics, as well as showing reconstructed epipolar images. See Figure B.1 and Table B.1.

As hypothesized, the models trained with the object plane closer to $z = 0$ perform better qualitatively and quantitatively. While this is, perhaps, an obvious result, we believe that it is an important one. In particular, it means that light fields with worse initial parameterization are more difficult to learn, and supports the result that learning re-parameterization via local affine transforms vastly improves neural light field interpolation quality.

### B.2. Non-Axis-Aligned Positional Encoding

Positional encoding with non-axis-aligned interpolation kernels (e.g. Gaussian positional encoding [53]) could be seen as a potential way around the issues discussed above. However, it is important to note that positional encoding with *any fixed set* of interpolation kernels will break down for certain scenes with different depths/depth ranges. For example, this is the case when the interpolation kernels do not align with the light field's color level sets, or there are not enough interpolation kernels to represent high frequencies for particular directions in the light field. In other words, no single-set of interpolation kernels works for all scenes, and ray-space embedding effectively tunes the interpolation per-scene/per-region in ray space,

## C. Embedding Visualization

In Figure C.2 we visualize predicted views, predicted EPIs, and the embedded ray-space given as input to the color network for:

1. The baseline approach
2. Feature space embedding
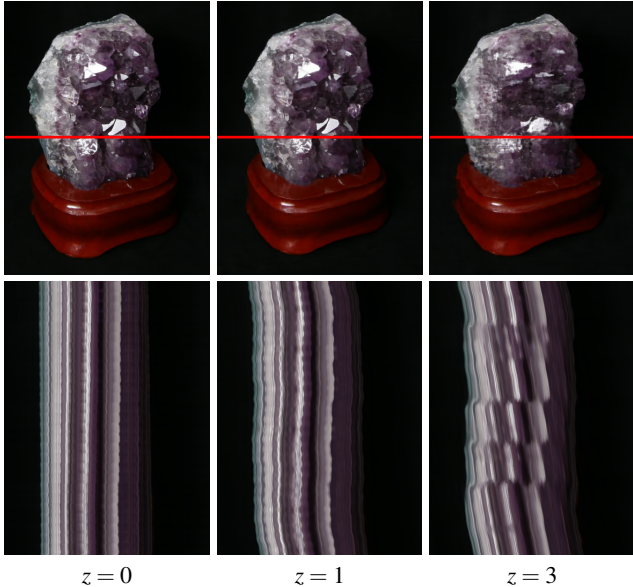
$z=0$       $z=1$       $z=3$

Figure B.1. **Effect of initial parameterization.** The *top* row shows predicted images, *bottom* predicted EPIs. Reconstruction becomes progressively worse for a worse initial parameterization ($\pi^{uv}$ moving further and further from the object plane at $z=0$).

3. Local affine transform embedding

For the embedding visualization, RGB colors denote the first three principal components of embedded ray-space.

Note the wiggling artifacts in the EPI predicted without embedding (*left*) . As discussed above, the baseline approach produces axis-aligned interpolation kernels which are ill-suited for interpolating slanted color level sets in the EPI.

The feature embedding network (*center*) learns what is essentially a set of texture coordinates for the object. It registers disparate rays that hit the same 3D points, and interpolates views fairly well, but does not guarantee multi-view consistency as the embedding network is still under-constrained for unobserved views.

Although the transforms themselves are not visualized (it is ray-space *after* applying the transforms that is visualized), the local affine transform network (*right*) predicts a set of transforms that are *largely constant*. This is because the depth range of the scene is limited, and a small set of re-parameterizations works well for all of ray-space. As the network learns a simpler output signal, it naturally interpolates more effectively, even to unobserved views. While the results for feature embedding and local affine transform embedding look similar in Figure C.2, we encourage readers to visit our web-page. It is far easier to see artifacts of the feature embedding approach in video form.
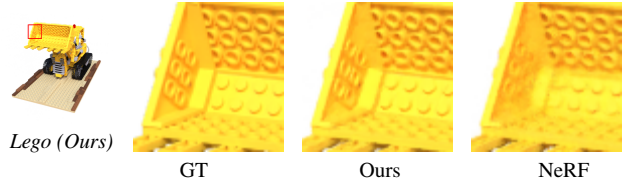


*Lego (Ours)*

GT      Ours      NeRF

Figure C.1. We achieve 29.14 dB on the *Lego* sequence from the NeRF Synthetic dataset [31] with a $32^3$ voxel grid at $800 \times 800$ pixel resolution, compared to 27.26 dB for AutoInt [23] with 32 sections, and 32.54 dB for NeRF.



No Embedding    Feature Embedding    Local Affine Embedding

Figure C.2. **Effect of different embeddings**. The *top* row shows predicted images, *middle* predicted EPIs, and *bottom* a visualization of the embedding space for each method.

## D. More Experimental Results

We show per-scene metrics in Tables D.1, D.2, D.3, and D.4. Additionally, we highly encourage readers to visit our project webpage, which contains image comparisons for every scene and video comparisons for a select few scenes.

**360° Scenes.** In Figure C.1 we show preliminary results for our $32^3$ subdivided model with Plücker parameterization applied to the *Lego* scene in the NeRF Synthetic [31] dataset. We use the same evaluation protocol as in NeRF for this scene. Our PSNR is slightly worse than NeRF, but better than AutoInt for the same grid resolution. Additionally, in some regions, we are able to better recover fine-grained texture on the *Lego* model.

**Student-Teacher Training.** We additionally provide results in Tables D.1 for our method when the input data is augmented with a 10x10 grid of renderings from a fully trained NeRF. We label this method as "Ours (w/t)," or our method with "student-teacher" training. With this approach, our method outperforms NeRF quantitatively in terms of PSNR, but at the cost of increased training time.

Table D.1. Per-scene breakdown results from NeRF's Real Forward-Facing dataset [31].

| Scene | PSNR↑ | | | | SSIM↑ | | | | LPIPS↓ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NeRF [31] | AutoInt [23] | Ours | Ours (w/ t) | NeRF [31] | AutoInt [23] | Ours | Ours (w/ t) | NeRF [31] | AutoInt [23] | Ours | Ours (w/ t) |
| Fern | **26.92** | 23.51 | 24.25 | <u>26.06</u> | **0.903** | 0.810 | 0.850 | <u>0.893</u> | **0.085** | 0.277 | 0.114 | <u>0.104</u> |
| Flower | 28.57 | 28.11 | <u>28.71</u> | **28.90** | 0.931 | 0.917 | **0.934** | **0.934** | 0.057 | 0.075 | **0.038** | <u>0.053</u> |
| Fortress | **32.94** | 28.95 | 31.46 | <u>32.60</u> | **0.962** | 0.910 | 0.954 | <u>0.961</u> | **0.024** | 0.107 | <u>0.027</u> | 0.028 |
| Horns | 29.26 | 27.64 | **30.12** | <u>29.76</u> | 0.947 | 0.908 | **0.955** | <u>0.952</u> | <u>0.058</u> | 0.177 | **0.044** | 0.062 |
| Leaves | **22.50** | 20.84 | 21.82 | <u>22.27</u> | <u>0.851</u> | 0.795 | 0.847 | **0.855** | <u>0.103</u> | 0.156 | **0.086** | 0.104 |
| Orchids | **21.37** | 17.30 | 20.29 | <u>21.10</u> | **0.800** | 0.583 | 0.766 | <u>0.794</u> | <u>0.108</u> | 0.302 | **0.103** | 0.113 |
| Room | <u>33.60</u> | 30.72 | 33.57 | **34.04** | <u>0.980</u> | 0.966 | 0.979 | **0.981** | 0.038 | 0.075 | **0.037** | **0.037** |
| T-rex | 28.26 | 27.18 | **29.41** | <u>28.80</u> | 0.953 | 0.931 | **0.959** | **0.959** | 0.049 | 0.080 | **0.034** | <u>0.040</u> |

Table D.2. Per-scene breakdown results from the *Undistorted* Real Forward-Facing dataset used in NeX [61]

| Scene | PSNR↑ | | SSIM↑ | | LPIPS↓ | |
|---|---|---|---|---|---|---|
| | NeX [61] | Ours | NeX [61] | Ours | NeX [61] | Ours |
| Fern | **26.46** | 24.49 | **0.913** | 0.856 | **0.068** | 0.107 |
| Flower | **29.39** | 28.93 | **0.947** | 0.937 | 0.041 | **0.033** |
| Fortress | **32.31** | 31.32 | **0.963** | 0.955 | **0.024** | 0.026 |
| Horns | 29.81 | **29.88** | **0.959** | 0.952 | **0.039** | 0.050 |
| Leaves | **22.66** | 21.62 | **0.879** | 0.845 | **0.082** | 0.082 |
| Orchids | **20.51** | 19.93 | **0.792** | 0.754 | **0.096** | 0.109 |
| Room | **33.40** | 33.24 | **0.979** | 0.978 | **0.033** | 0.036 |
| T-rex | 29.36 | **29.44** | **0.965** | 0.963 | 0.037 | **0.032** |

Table D.3. Per-scene breakdown results from NeX's Shiny Dataset [61]

| Scene | PSNR↑ | | SSIM↑ | | LPIPS↓ | |
|---|---|---|---|---|---|---|
| | NeX [61] | Ours | NeX [61] | Ours | NeX [61] | Ours |
| CD | 31.92 | **35.44** | 0.971 | **0.980** | **0.028** | **0.014** |
| Crest | **24.78** | 24.48 | **0.870** | 0.858 | **0.051** | 0.052 |
| Food | **25.61** | 25.21 | **0.905** | 0.885 | **0.048** | 0.053 |
| Giants | **28.50** | 27.99 | **0.946** | 0.930 | **0.038** | 0.039 |
| Lab | 31.20 | **34.39** | 0.965 | **0.982** | 0.031 | **0.013** |
| Pasta | **23.21** | 22.11 | **0.915** | 0.890 | **0.045** | 0.065 |
| Seasoning | **31.07** | 29.48 | **0.970** | 0.957 | **0.028** | 0.045 |
| Tools | **29.86** | 28.90 | **0.974** | 0.968 | **0.018** | 0.022 |

Table D.4. Per-scene breakdown results from the Stanford Light Field dataset [60]

| Scene | PSNR↑ | | | | SSIM↑ | | | | LPIPS↓ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NeRF [31] | X-Fields [2] | NeX [61] | Ours | NeRF [31] | X-Fields [2] | NeX [61] | Ours | NeRF [31] | X-Fields [2] | NeX [61] | Ours |
| Amethyst | <u>39.746</u> | 37.232 | 39.062 | **40.120** | 0.984 | 0.982 | 0.983 | **0.985** | 0.026 | 0.032 | <u>0.023</u> | **0.019** |
| Beans | **42.519** | 40.911 | 41.776 | <u>41.659</u> | **0.9944** | 0.9931 | <u>0.9938</u> | 0.9933 | <u>0.014</u> | 0.017 | 0.016 | **0.012** |
| Bracelet | <u>36.461</u> | 34.112 | 34.888 | **36.586** | <u>0.9909</u> | 0.9857 | 0.988 | **0.9913** | <u>0.0094</u> | 0.0260 | 0.0152 | **0.0087** |
| Bulldozer | <u>38.968</u> | 37.350 | 38.131 | **39.389** | 0.983 | <u>0.986</u> | 0.985 | **0.987** | 0.063 | 0.032 | <u>0.027</u> | **0.024** |
| Bunny | <u>43.370</u> | 42.251 | 42.722 | **43.591** | 0.9892 | <u>0.9894</u> | 0.9885 | **0.9900** | 0.029 | <u>0.022</u> | 0.036 | **0.013** |
| Chess | **41.146** | 37.996 | 39.938 | <u>40.910</u> | <u>0.9915</u> | 0.9882 | 0.9910 | **0.9920** | 0.028 | 0.034 | <u>0.020</u> | **0.016** |
| Flowers | <u>37.910</u> | 37.590 | 36.982 | **39.951** | 0.977 | <u>0.981</u> | 0.978 | **0.984** | 0.076 | <u>0.035</u> | 0.036 | **0.030** |
| Knights | **35.978** | 31.491 | <u>35.678</u> | 34.591 | **0.986** | 0.974 | <u>0.986</u> | 0.982 | **0.0142** | 0.0501 | 0.0168 | <u>0.0143</u> |
| Tarot (Small) | <u>34.221</u> | 30.830 | 33.134 | **36.046** | <u>0.982</u> | 0.975 | 0.981 | **0.989** | <u>0.014</u> | 0.033 | 0.016 | **0.006** |
| Tarot (Large) | **24.907** | 24.154 | 22.487 | <u>24.904</u> | <u>0.910</u> | 0.893 | 0.833 | **0.914** | <u>0.059</u> | 0.074 | 0.117 | **0.039** |
| Treasure | <u>34.761</u> | 33.904 | 32.350 | **37.465** | 0.972 | <u>0.977</u> | 0.967 | **0.982** | <u>0.027</u> | 0.041 | 0.040 | **0.019** |
| Truck | <u>40.723</u> | 38.883 | 38.292 | **41.440** | 0.986 | 0.984 | 0.986 | **0.989** | 0.087 | 0.042 | <u>0.038</u> | **0.033** |