# G-NeLF: Memory- and Data-Efficient Hybrid Neural Light Field for Novel View Synthesis
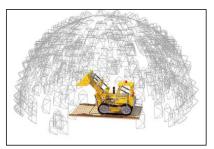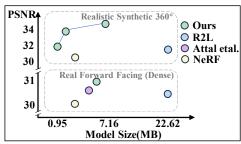
Lutao Jiang[1]    Lin Wang[1,2] *

[1] VLIS LAB, AI Thrust, HKUST(GZ)    [2]Dept. of CSE, HKUST

ljiang553@connect.hkust-gz.edu.cn, linwang@ust.hk

Project Page: https://vlislab22.github.io/G-NeLF/

**a) R2L: 10,000 Synthesis Views**    **b) Ours: Only 100 Original Views**    **c) Comparison of PSNR and Model Size**

Figure 1. **A comparison between R2L [40] and our G-NeLF**. a) R2L needs to use a pre-trained NeRF model to synthesize 10,000 training images to train the NeLF model. b) Our G-NeLF only needs the original amount of images the same as NeRF's dataset. c) The comparison of rendering quality and model size among ours, R2L, and NeLF proposed by Attal *et al*. [1]. The upper half is the Realistic Synthetic 360° dataset. The lower half is scenes containing relatively dense views in the Real Forward-Facing Dataset.

## Abstract

*Following the burgeoning interest in implicit neural representation, Neural Light Field (NeLF) has been introduced to predict the color of a ray directly. Unlike Neural Radiance Field (NeRF), NeLF does not create a point-wise representation by predicting color and volume density for each point in space. However, the current NeLF methods face a challenge as they need to train a NeRF model first and then synthesize over 10K views to train NeLF for improved performance. Additionally, the rendering quality of NeLF methods is lower compared to NeRF methods. In this paper, we propose **G-NeLF**, a versatile grid-based NeLF approach that utilizes spatial-aware features to unleash the potential of the neural network's **inference capability**, and consequently overcome the difficulties of NeLF training. Specifically, we employ a spatial-aware feature sequence derived from a meticulously crafted grid as the ray's representation. Drawing from our empirical studies on the adaptability of multi-resolution hash tables, we introduce a novel grid-based ray representation for NeLF that can represent the entire space with a very limited number of parameters. To better utilize the sequence feature, we design a lightweight ray color decoder that simulates the ray propagation process, enabling a more efficient inference of the ray's color. G-NeLF can be trained **without necessitating significant storage overhead** and with **the model size of only 0.95 MB** to surpass previous state-of-the-art NeLF. Moreover, compared with grid-based NeRF methods, e.g., Instant-NGP, we only utilize one-tenth of its parameters to achieve higher performance. Our code will be released upon acceptance.*

## 1. Introduction

Driven by inspiration from Implicit Neural Representation (INR) [27, 30], the task of novel view synthesis – representing a 3D scene from 2D views – has recently witnessed a significant breakthrough with the introduction of Neural Radiance Fields (NeRF) [3, 7, 8, 13, 14, 18, 23, 26, 28, 29, 31, 43–45, 49, 50] and Neural Light Field (NeLF) [1, 5, 16, 32, 33, 40, 47]. *A key distinction between NeRF and NeLF is the necessity for the physical modeling of each point.* NeRF requires querying the color and volume density of each point in the space from a neural network and then utilizing the classical volume rendering [20] to render an image. In contrast, the NeLF network predicts the ray's color directly based on the ray coordinate, *e.g.*, the coordinates of two points can encode a ray's coordinate. However, since there are no constraints based on physical world modeling, it is difficult for the neural network to learn the color of each ray.

---

*Corresponding author.

A seminal work, R2L [40], proposes a deep MLP network to represent the light field faithfully by leveraging the network's greater *memory capability* and expressivity. Unfortunately, this causes difficulties in training such a cumbersome network that only uses 2D images with 100 views (commonly used for training NeRF). For this reason, it synthesizes 10K images with various viewpoints using a pre-trained NeRF. The following NeLF methods [5, 16] also adhere to this design. However, such an image synthesis operation takes a long time and consumes storage for the record rays' coordinates and color. This renders challenges for applications that require images with higher resolutions. Consequently, *it remains challenges of how to reduce computation cost and alleviate the training difficulty of NeLF under the condition of data shortage*.

In this paper, we present **G-NeLF**, a lightweight NeLF approach hybridized with spatial-aware grids, which overcomes the training difficulties of R2L with only 100 views or fewer for training, as demonstrated in the comparison in Fig.1 (a) and (b). Also, G-NeLF can be trained without requiring significant storage overhead with the model size of only **0.95 MB**, see Fig. 1(c). *Our key idea* is to utilize spatial-aware feature sequence as the ray representation method to unleash the potential of the recurrent network's *inference capability* by leveraging the *temporal order information* in ray propagation.

To obtain the feature sequence, we design a 3D spatial-aware feature grid in the space. Instant-NGP [29] introduces a multi-resolution hash table designed to store features of 3D grids. While the hash table reduces the parameter count compared with the vanilla 3D grid, its quantity remains relatively large. Our analysis and empirical evidence (See Tab. 1 and Fig. 4) reveal that in Instant-NGP's hash table, low-resolution features primarily encode shape information, while high-resolution features predominantly capture details related to surface appearance and color. We analyze it in Sec. 3.1. Though we don't model the point-wise volume density representation, we can still assimilate this property into our architecture by the integration of multiple hash-based multi-resolution grids, benefiting from the *spatial awareness of our ray representation*. Based on this, we present a 3D grid decomposition representation for NeLF, which only comprises a compact set of 474,432 parameters. Given that ray propagation inherently possesses strong temporal order information, efficiently capitalizing on this aspect is crucial to achieving our objectives. Therefore, we design an LSTM-based [17] ray color decoder (Sec. 3.2), a structure that experiments suggested is optimal, to predict the ray's color throughout the sequence.

Our contribution can be summarized as follows:

- We present G-NeLF, an *implicit* approach for the novel view synthesis task, *i.e.*, not depend on intricate point-wise volume-based physical modeling techniques com-

monly used in NeRFs, and *not depend on other auxiliary data [40] for training*.
- We propose a novel grid-based ray representation method, which can unleash the potential of *inference ability* of the neural networks. Also, we design a lightweight ray color decoder, which simulates the temporal nature of ray propagation and is extremely *flexible* to adjust the model size.
- Drawing on the analysis of the integration between multi-resolution and hash mapping, we propose a novel 3D grid representation for NeLF characterized by an *exceptionally small number of parameters*.
- We conduct extensive experiments to prove the superiority of our G-NeLF. As depicted in Fig.1 (c), compared with SoTA NeLF methods, *e.g.*, R2L [40], we surpass it by **0.04dB** only with **100** views and **0.95** MB model size, and surpass it by **2.48dB** with the increase of model size; Compared with SoTA lightweight NeRF compression method, VQRF [23], we use same model size, **1.4MB**, to surpass it by **0.79dB**; Compared with NeRF methods [3, 7, 29], we achieve comparable results.

## 2. Related Works

**Neural Light Field.** The concept of the light field was initially introduced by [15, 22]. Both studies employ a 4-D coordinate system $((u, v), (s, t))$ to represent a ray within space. By pinpointing the positions of $(u, v)$ and $(s, t)$ on two distinct planes, the ray is consequently determined. The recent advancements in NeRF [28] have catalyzed interest in Implicit Neural Representation (INR). Consequently, an increasing number of researchers are recognizing the potential of INR in the domain of the light field. Light Field Networks (LFNs) [32] pioneers the integration of INR with the light field. They utilize meta-learning to get the prior of multi-view consistent light field and use one single view image to reconstruct the other views. But compared with NeRFs, it remains a relatively low quality. Suhail *et al*. [33] utilize transformer [38] structure to factor in epipolar lines from various views for novel views. But its inference speed and training cost is unacceptable. Attal *et al*. [1] propose to map 4D ray to an embedding space. They use some voxels to represent the local light field while simultaneously learning the opacity of the voxels. R2L [40] introduces the innovative method of using $K$ points, as opposed to the conventional two points, to represent a ray. To mitigate the data constraints inherent in NeLF training, they synthesize a lot of pseudo data from the pre-trained NeRF model. DyLiN [47] first makes the light field dynamic depending on a deformation network. MobileR2L [5] incorporates a $1 \times 1$ convolution layer as a substitute for the MLP used in R2L, and integrates multiple lightweight super-resolution modules to further reduce the computation consumption, which makes it run on iPhone. LightSpeed [16] adopts the K-Planes [11] approach, employing six feature planes to encode the 4-D
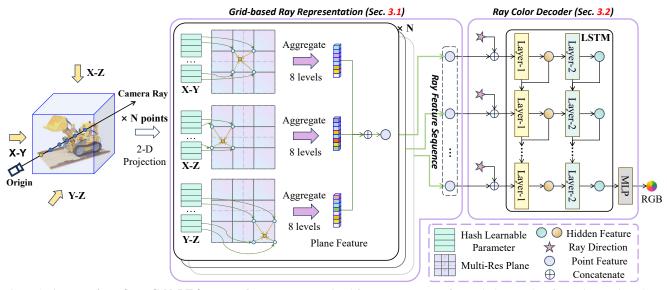
Figure 2. **An overview of our G-NeLF framework**. For one ray emitted from a camera, we first orderly sample a few points on it and use our designed hash multi-resolution triplane to obtain their feature. Then these features will compose the ray's feature sequence to represent this ray. Finally, we design a ray color decoder to transform the ray feature sequence into RGB color.

two-plane coordinates of rays, while simultaneously integrating a super-resolution module [5] to increase the resolution. However, compared with R2L and MobileR2L, LightSpeed imposes an increasing demand for the training dataset of $30K$ pseudo images.

Differently, we harness spatial information to unleash the potential of the inference ability of neural networks. This solves the problems of previous works, including training data requirements and the complexity of the whole training pipeline. Moreover, we only use the original training data with a 0.95 MB model size to surpass R2L.

**Grid Design in NeRF.** NeRF [28] is proposed for the novel view synthesis task and 3D reconstruction task, and it has become the dominant paradigm. However, the low training and inference speed present challenges for real-world applications. Recently, a lot of methods are proposed to solve these problems mainly from point sampling aspect [19, 24, 25] and data structure aspect including OcTree [2, 41, 46], mesh [8, 37] and grid representation [4, 7, 11, 12, 14, 23, 29, 35, 36, 42]. Instant-NGP [29] represents the most advanced fusion of NeRF and explicit grid representation, which maps the multi-resolution grid to a table using hash function. Although the hash table has reduced the number of parameters than the naive 3D grid, the amount is still too large. EG3D [6] introduces the tri-plane to represent a NeRF Grid. TensoRF [7] proposes a 3D decomposition method that uses three planes and three axes to represent a 3D grid. However, due to its expansive feature dimension width, the amount of parameters is still too large. VQRF [23] uses vector quantization and weight quantization techniques to compress the Grid to 1MB. K-Planes [11] and HexPlane [4] extend the grid design to dynamic scenes by composing several 2D planes of X, Y, Z, and T axes.

Different from these methods, we delve deeper into the role of the hash mapping function during the training phase and demonstrate its adaptability and flexibility. Building on this foundation, we formulate a highly efficient grid representation for NeLF, requiring only about 0.47 million parameters to surpass the previous methods, *e.g.*, R2L. While NeRF methods typically demand more sampling points to ensure precise volume density modeling and spatial stability, our G-NeLF simplifies the process. Our G-NeLF only needs to sample fewer points' features to represent one ray's feature, avoiding more complex training requirements.

## 3. Method

**Overview.** The goal of G-NeLF is to synthesize arbitrary views of the scene. The design is guided by three key ideas. **1)** We employ a feature sequence as a substitute for the previous coordinate-based ray representation. **2)** Based on our analysis, we design a hash-based multi-resolution tri-plane. **3)** We rely on the strong temporal information embedded within the ray's feature sequence to infer its color.

Our workflow, as depicted in Fig. 2, encompasses the following steps. We emit rays based on the camera model to form an image, with each ray corresponding to a pixel. Initially, the ray crosses a bounding box in space. Next, we sequentially sample some points on the ray and project them onto three orthogonal feature planes. Utilizing the feature grid we designed, we obtain the feature vector for each point and systematically assemble these vectors to form the ray's feature sequence (Sec. 3.1). To decode the ray's color, we design a lightweight ray color decoder based on an LSTM
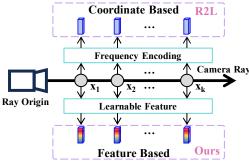
3

Figure 3. **Illustration of ray representation method difference** between R2L [40] and our G-NeLF. R2L directly concatenates frequency encoded [28] coordinates as the network input. We use feature sequence to represent one ray.

network (Sec. 3.2). This decoder receives the ray's feature sequence along with the ray's direction as inputs and predicts the ray's color.

## 3.1. Grid-based Ray Representation

**Ray Representation.** Early ray representation methods [15, 22, 32, 33] primarily focus on 4-D coordinate $((u, v), (s, t))$ representation. Recent approaches [5, 40, 47] begin to employ the $3K$-D coordinate representation method. These method samples $K$ points, $\{\mathbf{x_1}, ..., \mathbf{x_k}\}$, along the ray to construct the representation, as illustrated in the upper part of Fig. 3. However, reliance solely on coordinate representation poses challenges in learning the mapping relationship between a ray's color and its coordinates, complicating the accurate prediction of color. Consequently, these methods require extensive training data to effectively leverage the neural network's memory capacity for encoding the color of rays.

Rethinking the ignorance of utilizing the neural networks' ***inference capability*** [5, 16, 40, 47], we aim to develop a representation method that effectively addresses this aspect. We introduce a $KN$-D representation, wherein ***a sequence of features is employed to represent a single ray***, as illustrated in the below of Fig. 3. Specifically, we orderly sample $K$ points along the ray, and extract their corresponding $N$-dimensional spatial-aware features to compose a feature sequence, $\{\mathcal{F_1}, ..., \mathcal{F_k}\}$. The spatial awareness embedded within this feature sequence is key to enabling effective inference. Moreover, although there is no explicitly pointwise physical modeling, the spatial-aware feature sequence still retains the capability to constrain the view consistency.

**Feature Grid Design.** To derive the feature sequence, $\{\mathcal{F_1}, ..., \mathcal{F_k}\}$, we develop a notably lightweight grid decomposition method, which only requires about 0.47 M parameters. To guarantee that each point along the ray is endowed with distinct learnable features, we configure a conceptual grid in the space. Each vertex of the grid holds a learnable feature vector. We then proceed to compute the weighted average of features for each point based on its distance to

| Mask | Cosine Similarity ↓ | PSNR ↑ |
|---|---|---|
| w/o mask | - | 34.64 |
| Top-6 | 0.9739 | 28.08 |
| Top-10 | 0.8234 | 17.00 |
| w/o mask | - | 28.20 |
| Top-6 | 0.9860 | 25.74 |
| Top-9 | 0.8287 | 17.65 |

Table 1. **Comparison of different numbers of feature masks.** *w/o mask* means original Instant-NGP model. *Top-6* means that the top 6 resolutions are masked to 0. *Top-10* means that the top 10 resolutions are masked to 0.
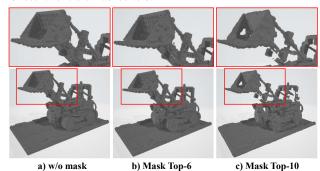


**a) w/o mask**     **b) Mask Top-6**     **c) Mask Top-10**

Figure 4. **Visual comparison of different numbers of feature masks.** a) Instant-NGP. b) Mask top 6 resolution in Instant-NGP. c) Mask top 10 resolution in Instant-NGP. Observing the results, it is evident that masking features at higher resolutions still allows for the preservation of overall shape information.

the vertices of the grid where it is situated, endowing it with a unique feature vector. Based on our analysis, we design a hash multi-resolution tri-plane as the decomposition method for our conceptual grid.

We conduct a thorough analysis of the relationship between multi-resolution and hash mapping based on Instant-NGP [29], exploring the mechanisms underlying their interaction. We argue that ***the low-resolution primarily determines the shape features while the high-resolution represents the surface appearance features, self-adapting at the training stage gradually***. In the hash table, the low-resolution grids are exclusively owned by themselves. Therefore, they can reserve the structure information accurately. On the contrary, in the hash table of the highest resolution grid, $2^{19}$ items are used to represent $2^{36}$ positions' feature, which means approximately $100,000$ positions are mapped to a single item in the hash table. Because the positions far away from the first surface don't influence the color in the volume rendering formula [20, 28], they don't affect the corresponding color feature entries in the hash table during the training phase. As a result, in the set of positions that are mapped to the same entry in the hash table, the position carrying the larger weight becomes the dominant one in influencing this entry following training progresses. Consequently, the features at higher resolutions adaptively concentrate on capturing the appearance features. To prove our argument, we conduct well-designed experiments.

To obtain the main function of different resolutions, we employ zero-masking on high-resolution features during the computation process. As demonstrated in Fig. 4, the shape is largely preserved even when the top six resolutions are masked. As shown in the top half of Tab. 1, the similarity between the original and masked hash grids is 0.9739. However, the PSNR metric, indicative of rendered image similarity, exhibits a significant reduction. When we mask the top ten resolutions, the shape is still relatively well-preserved with a similarity of 0.8234, but there is a more pronounced loss in color information. A similar trend can be observed when replacing the hash grid with our hash tri-plane in the bottom part of Tab. 1.

Though there is no explicit modeling of volume density in G-NeLF, the feature sequence is spatial-aware, which creates a huge potential to apply the same strategy. Even, without the requirements of extra information, *i.e.*, volume density, NeLF has a greater potential to allow a much more compact representation. For NeLF, we design a more compact hash-based multi-resolution tri-plane representation that yields superior results compared with the hash grid. As depicted in Fig. 2, our process begins by projecting the points from the ray onto three orthogonal planes. Subsequently, according to the projected coordinates, we identify the corresponding grid vertices across all resolutions for each feature plane. Then, we get the vertices' features of each resolution from a hash table. The hash table mapping function is

$$Index = (xa \oplus yb) \mod T \quad (1)$$

where $x$ and $y$ are the indexes of height and width, $a = 1$, $b = 2,654,435,761$, $\oplus$ is the XOR opration, $T$ is the hash table size. We use the $Index$ to retrieve the specific feature vector from the corresponding hash table. Subsequently, we calculate the weighted-sum feature based on the distance between the point and these vertices. Finally, we concatenate the feature vectors from all resolutions as the feature of that point. The concrete procedure of calculation is illustrated in Algorithm 1.

## 3.2. Ray Color Decoder

Once we acquire the feature sequence, $\{\mathcal{F}_1, ..., \mathcal{F}_k\}$, we can predict its color from this sequence. A direct approach is to employ an MLP network to directly predict the color. However, due to the neglect of the potent temporal-order sequential information remaining in the feature sequence, the MLP network is not efficient in predicting the color. We also attempt to employ a positional-aware Transformer-based [38] architecture. Nevertheless, as same as the analysis in [48], since the calculation of attention is permutation-invariant and anti-order, it is not suitable for our target either. Despite there being some position encoding techniques, there remains an inevitable loss of certain temporal-order sequential features.

---

**Algorithm 1** Point's Feature Acquirement

1: **Input**: $x$, the 3D coordinate of input point; $\{\pi_{xy}, \pi_{xz}, \pi_{yz}\}$, the tri-plane; $s_i$, the grid size of $i$-th resolution; $N$, the number of resolutions of each plane; $hash_{\pi,i}(\cdot)$, the hash table of $i$-th resolution in the plane $\pi$;

2: **Output**: $F$, feature vector of the input point $x$;

3: **for** $\pi$ in $\{\pi_{xy}, \pi_{xz}, \pi_{yz}\}$ **do**

4:     $F_\pi = []$;   # empty list

5:     **for** $i \leftarrow 1$ to $N$ **do**

6:         $x_\pi$ = Projection(x, $\pi$);   # 2D Coordinate

7:         $t\_l = hash_{\pi,i}([x_\pi[0] \mathbin{/\!/} s_i, x_\pi[1] \mathbin{/\!/} s_i])$;

8:         $t\_r = hash_{\pi,i}([x_\pi[0] \mathbin{/\!/} s_i, x_\pi[1] \mathbin{/\!/} s_i + 1])$;

9:         $b\_l = hash_{\pi,i}([x_\pi[0] \mathbin{/\!/} s_i + 1, x_\pi[1] \mathbin{/\!/} s_i])$;

10:        $b\_r = hash_{\pi,i}([x_\pi[0] \mathbin{/\!/} s_i + 1, x_\pi[1] \mathbin{/\!/} s_i + 1])$;

11:        $F_\pi$.append(Weighted_Sum($[t\_l, t\_r, b\_l, b\_r]$));

12:     **end for**

13:     $F_\pi$ = Concat($F_\pi$);

14: **end for**

15: $F$ = Concat($[F_{\pi_{xy}}, F_{\pi_{xz}}, F_{\pi_{yz}}]$);

16: **return** $F$;

17: **End**.

---

Considering that ray propagation in the real world follows a path from near to far, we propose employing a lightweight LSTM network [17] to emulate this sequential process, as shown in Fig. 2 (right). Unlike text sequences, the temporal order of a ray's propagation is crucial. The results presented in Tab. 7 demonstrate that recurrent computation is a critical component enabling the effectiveness of our ray color decoder. Andrej [21] *et al.* argue that LSTM with more than two layers exhibits fundamental differences from its single-layer counterparts. Consistent with this assertion, our experimental results (Fig. 7(a)) indicate that a 2-layer LSTM surpasses the performance of a 1-layer LSTM when computational costs are equivalent. To simulate the view-dependent phenomenon, we additionally incorporate the ray direction into the decoder. We employ Sphere Harmonics, as adopted by Ref-NeRF [39] and Instant-NGP [29], to encode the 3D normalized direction vector. The encoded direction is concatenated to the ray sequential feature. Formally, the decoding process can be written as

$$\mathcal{H}_i, \mathcal{C}_i = LSTM(Concat(\mathcal{F}_i, d), \mathcal{H}_{i-1}, \mathcal{C}_{i-1}), \quad (2)$$

$$(r, g, b) = MLP(\mathcal{H}_k), \quad (3)$$

where $i = 1, 2, ..., k$, $\mathcal{F}_i$ is coming from ray's feature sequence $\{\mathcal{F}_1, ..., \mathcal{F}_k\}$, $d$ is the encoded ray's direction, $\mathcal{H}_i$ and $\mathcal{C}_i$ are hidden state and cell state in LSTM, $LSTM(\cdot)$

is a LSTM network, $Concat(\cdot)$ is the concatenation function, $MLP(\cdot)$ is a two-layer MLP network, $(r, g, b)$ is the target, ray's color.

## 4. Experiments

### 4.1. Datasets and Implementation Details

**Datasets.** We use two datasets to evaluate, including Realistic Synthetic 360° and Real Forward-Facing [28]. The Realistic Synthetic 360° comprises 8 scenes. Each scene provides 100 training views, 100 validation views, and 200 test views, all of which have a resolution of $800 \times 800$. To compare with R2L [40], we downsample all images by 2x during training and testing. To further evaluate our method at high resolution, we use the original resolution to compare with NeRF methods. The Real Forward-Facing also consists of 8 scenes, with each scene containing between 20 to 62 images. As same as other methods, we hold every 8th image for the test. We downsample all images in Forward-Facing by a factor of 8.

**Implementation Details.** All our experiments were conducted on a single RTX3090 GPU using float16 accuracy. We provide three configurations for experiments of 360° scenes. For the small and medium model settings, we use a grid with 8 levels' resolution from 16 to 1024, the max hash table size is $2^{14}$, and the feature dimension of each level is 2. The ray color decoder is a two-layer LSTM network with 32 (small) and 128 (medium) hidden sizes. For the large version, we use 16 levels' resolution with 16 to 2048, and the max hash table size is $2^{16}$. The ray color decoder has three layers with 128 hidden sizes. For 360° scenes, the feature sequence length is set at 256. For forward-facing scenes, we convert the coordinate system to NDC and establish the feature sequence length at 128. We use the MSE color loss to supervise the whole model in an end-to-end manner.

### 4.2. Overall Comparison

**Realistic Synthetic 360° Dataset.** We first compare our method with the NeLF method, R2L [40], at $400 \times 400$. In Tab. 9, we show PSNR and SSIM comparison results scene by scene. Different from the training difficulty and resource consumption that R2L needs, we only need the original training data, such as 100 views, saving over 56 GB space. This advantage also enables our method to train models on higher-resolution data effectively. Despite our method's minimal resource consumption, it significantly outperforms R2L, underscoring its superiority. Our small configuration, utilizing just 100 views for training, still surpasses R2L 0.04 dB with only a 0.95 MB model size. Furthermore, continuing to scale up the model size results in a consistent improvement, whereas the performance of R2L has reached its saturation point. In Fig. 5, we present a visual comparative analysis between R2L and Ours. For example, in the *Lego* scene, R2L encounters difficulties rendering the region behind the shovel, producing an image
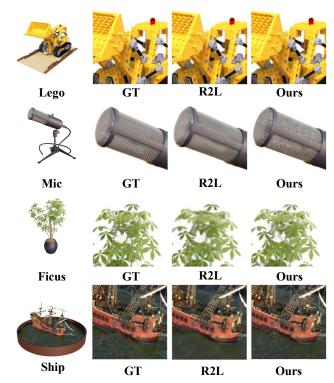


Figure 5. **Visual comparison** with R2L [40] on *Lego*, *Mic*, *Ficus*, and *Ship* scene in Realistic Synthetic 360° dataset.

with compromised clarity. Our method, conversely, renders this area with distinct sharpness, showcasing its enhanced performance.

As our G-NeLF is a lightweight model, we also compare it with the SoTA NeRF model compression method, VQRF [23]. As Tab. 10 shown, when evaluated under comparable model size, 1.43 MB for *DVGO-VQDF* and 1.41 MB for *Ours-M*, our method consistently delivers superior results by **+0.79** PSNR. This large gain shows the efficiency of our tiny hash tri-plane design and ray color decoder.

To demonstrate that we have pushed the NeLF far beyond the performance of previous methods, we also compare it with some SoTA grid-based NeRF methods, including TensoRF [7] and Instant-NGP [29]. As shown in Tab. 10, our method can achieve comparable results on par with SoTA NeRF techniques, which proves that NeLF also has a large potential on novel view synthesis tasks with relatively low resources. Grid-based NeRFs usually introduce a lot of storage consumption of model size, *e.g.*, 64.1MB of Instant-NGP [29] and 71.8MB of TensoRF [7]. Our ray-wise representation, which does not rely on precise physical modeling, enables us to represent a scene with a much more compact structure, about one-tenth parameter usage.

**Real Forward-Facing Dataset.** As shown in Tab. 12, we first compare our method with NeRF [28], NeLF [1], and R2L [40] on the Real Forward-Facing dataset at $504 \times 378$ resolution. At the head of the table, the number of original training views is in the brackets. Moreover, the results

| Method | Avg | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship | Size(MB) | Data Size |
|---|---|---|---|---|---|---|---|---|---|---|---|
| NeRF [28] | 30.47 | 33.90 | 25.56 | 28.88 | 34.64 | 31.42 | 29.22 | 30.84 | 29.30 | 4.56 | 100 |
| R2L [40] | 31.87 | 36.71 | 26.03 | 28.63 | 38.07 | 32.53 | 30.20 | 32.80 | 29.98 | 22.62 | 10,100 |
| Ours-S | 31.91 | 34.39 | 24.81 | 29.40 | 37.62 | 33.83 | 30.54 | 34.78 | 29.93 | 0.95 | 100 |
| Ours-M | 33.89 | 36.65 | 25.86 | 30.82 | 39.40 | 36.53 | 33.54 | 37.05 | 31.30 | 1.41 | 100 |
| Ours-L | 34.35 | 36.64 | 25.91 | 32.13 | 39.92 | 37.56 | 33.55 | 37.08 | 32.01 | 7.16 | 100 |

Table 2. PSNR quantitative comparison at $400 \times 400$ resolution on Realistic Synthetic 360° dataset. *Size* is model size. *Data Size* is the number of the training images. value means first, value means second, value means third.

| Category | Method | Avg | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship | Size(MB) | Data Size |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Compression | VQRF [23] | 31.77 | 33.80 | 25.38 | 32.67 | 36.47 | 34.27 | 29.28 | 33.11 | 29.24 | 1.43 | 100 |
| | **Ours-M** | 32.56 | 33.41 | 24.91 | 31.78 | 37.72 | 35.41 | 32.25 | 35.73 | 29.24 | **1.41** | 100 |
| NeRF | NeRF [28] | 31.01 | 33.00 | 25.01 | 30.13 | 36.18 | 32.54 | 29.62 | 32.91 | 28.65 | 5 | 100 |
| | Mip-NeRF [3] | 33.09 | 35.14 | 25.48 | 33.29 | 37.48 | 35.70 | 30.71 | 36.51 | 30.41 | 2.5 | 100 |
| | TensoRF [7] | 33.14 | 35.76 | 26.01 | 33.99 | 37.41 | 36.46 | 30.12 | 34.61 | 30.77 | 71.8 | 100 |
| | I-NGP [29] | 33.18 | 35.00 | 26.02 | 33.51 | 37.40 | 36.39 | 29.78 | 36.22 | 31.10 | 64.1 | 100 |
| NeLF | MobileR2L [5] | 31.34 | 33.66 | 25.05 | 29.80 | 36.84 | 32.18 | 30.54 | 34.37 | 28.75 | 8.20 | 10,100 |
| | LightSpeed [16] | 32.23 | 34.21 | 25.63 | 32.82 | 36.77 | 34.35 | 29.51 | 35.65 | 28.90 | 13.60 | 30,100 |
| | **Ours-L** | 33.19 | 34.56 | 25.01 | 32.45 | 38.03 | 36.78 | 32.32 | 36.23 | 30.10 | 7.16 | 100 |

Table 3. PSNR quantitative comparison at $800 \times 800$ resolution on Realistic Synthetic 360° dataset. The VQRF is based on DVGO [34]. *Size* is model size. *Data Size* is the number of training images. **The first part** is NeRF's compression method. **The second part** is the SoTA NeRF methods. **The third part** is the SoTA NeLF methods.

| Method | Dense | | | | | | Sparse | | | Additional |
|---|---|---|---|---|---|---|---|---|---|---|
| | Avg | Horns(54) | Trex(48) | Room(36) | Fortress(36) | Flower(29) | Leaves(21) | Orchids(21) | Fern(17) | Images |
| NeRF [28] | 30.17 | 28.86 | 28.10 | 33.07 | 32.61 | 28.22 | 22.40 | 21.29 | 26.86 | 0 |
| Attal *et al.* [1] | 30.65 | 30.12 | 29.41 | 33.57 | 31.46 | 28.71 | 21.82 | 20.29 | 24.25 | 0 |
| R2L [40] | 30.35 | 28.95 | 28.12 | 33.30 | 32.71 | 28.67 | 22.71 | 21.46 | 26.87 | 10,000 |
| Ours | 30.97 | 30.55 | 29.30 | 33.28 | 32.46 | 29.27 | 22.04 | 19.46 | 25.29 | 0 |

Table 4. PSNR quantitative comparison at $504 \times 378$ resolution on Real Forward-Facing. *Additional Images* is the required number of extra synthesized images. The number of original training views is in the brackets.

| Category | Method | Dense | | | | | | Sparse | | | Additional |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Avg | Horns(54) | Trex(48) | Room(36) | Fortress(36) | Flower(29) | Leaves(21) | Orchids(21) | Fern(17) | Images |
| NeRF | NeRF [28] | 29.10 | 27.45 | 26.80 | 32.70 | 31.16 | 27.40 | 20.92 | 20.36 | 25.17 | 0 |
| | Plenoxels [10] | 28.64 | 27.58 | 26.48 | 30.22 | 31.09 | 27.83 | 21.41 | 20.24 | 25.46 | 0 |
| | TensoRF [7] | 29.08 | 27.64 | 26.61 | 31.80 | 31.14 | 28.22 | 21.34 | 20.02 | 25.31 | 0 |
| NeLF | MobileR2L [5] | 28.85 | 27.01 | 26.71 | 32.09 | 30.81 | 27.61 | 20.52 | 20.06 | 24.39 | 10,000 |
| | LightSpeed [16] | 29.12 | 27.04 | 26.93 | 32.32 | 31.45 | 27.88 | 21.01 | 20.33 | 25.05 | 10,000 |
| | Ours | 29.26 | 28.08 | 27.20 | 31.89 | 30.77 | 28.35 | 20.88 | 19.01 | 23.68 | 0 |

Table 5. PSNR quantitative comparison at $1008 \times 756$ resolution on Real Forward-Facing. *Additional Images* is the required number of extra synthesized images. The number of original training views is in the brackets.

at $1008 \times 756$ resolution are shown in Tab. 5, including NeRF-based methods [7, 10, 23, 28] and NeLF-based methods [5, 16]. Unlike the meticulously constructed Realistic Synthetic 360° dataset, the Real Forward-Facing dataset, sourced from the real world, has inherent inaccuracies in camera parameters. This dataset also presents a higher scene complexity. We attribute the superiority over other methods to the fact that our technique does not hinge on point-wise modeling, allowing for some leniency in camera parameter estimation. Our method is fundamentally rooted in ray-level representations, which reduces the emphasis on exact camera parameter accuracy. In Fig. 6, we show the visual comparison with R2L.

## 4.3. Ablation Studies

**Design of Feature Grid.** To prove the superiority of our hash tri-plane design for NeLF, we replace it with the grid designs in Instant-NGP [29] and TensoRF [7] for comparison. All experiments are based on our 3-layer 128-width ray color decoder, which ensures ample capacity to ensure that the decoder does not become a bottleneck. The feature sequence length is set as 256. As Tab. 6 shows, our method still outperforms the two methods with such a small number of parameters. We believe this is due to more compact representations and easier training with fewer parameters.

**Design of Ray Color Decoder.** In Tab. 7, we present a comparison among different designs of ray color decoders,
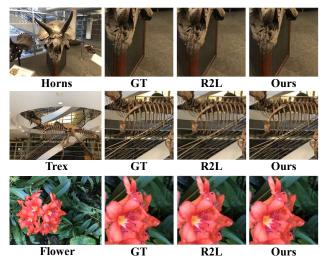
**Figure 6. Visual comparison** with R2L [40] on *Horns*, *Trex*, and *Flower* scene in Real Forward-Facing dataset.

| Method | PSNR | Params |
|---|---|---|
| Instant-NGP | 34.83 | 12,239,728 |
| TensoRF | 34.85 | 13,003,200 |
| Ours | **35.36** | **474,432** |

**Table 6. The grid design comparison** among Instant-NGP-based [29], TensoRF-based [7], and our hash tri-plane for our G-NeLF.

| Method | PSNR | Params | FPS |
|---|---|---|---|
| MLP-Based | 22.35 | 71,571,459 | 4.55 |
| RNN-Based | 15.52 | 25,347 | 6.54 |
| Transformer-Based | 31.56 | 21,603 | 1.32 |
| GRU-Based | 33.07 | 18,051 | 5.71 |
| Ours-LSTM Design | **33.82** | 23,299 | **6.71** |

Table 7. **The validation** of our ray color decoder design.

| Method | PSNR |
|---|---|
| w/o direction | 33.45 |
| Frequency Encoding | 33.40(-0.05) |
| Spherical Harmanics | 33.82(+0.37) |

Table 8. **The different encoding methods** for direction vector.

ensuring a similar computational load across all methods for fairness. Our experiments span many prevalent network architectures, including MLP, RNN, GRU [9], Transformer [38], and LSTM [17]. Our findings reveal that GRU and LSTM-based ray color decoders are particularly effective, validating our hypothesis that temporal sequential information is the key factor in inferring the ray feature sequence effectively. Moreover, although Transformer-based decoders incorporate position encoding techniques to perceive temporal relationships, the performance is still unsatisfactory for a lot of temporal prediction tasks [48].

**The Effect of Ray Direction.** To ascertain the influence of ray direction on the performance of the ray color decoder, we conduct the ablation experiments under our G-NeLF small model. We compare the different direction vector encoding methods, including Frequency Encoding
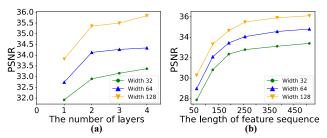


Figure 7. (a) The ablation study of the number of layers. (b) The ablation study of the length of feature sequence.

used in vanilla NeRF [28] and Spherical Harmonics used in Instant-NGP [29]. As shown in Tab. 8, Spherical Harmonics is more suitable for our G-NeLF.

**The Number of LSTM Layers.** According to the research of Andrej *et al*. [21], the LSTM network with more than two layers exhibits fundamental differences and significant performance enhancements compared with its single-layer counterparts. We use the small configuration grid setting and we set the sequence length as 256. Fig. 7 (a) presents a comparative analysis of the impact of varying the number of layers across three different width settings: 32, 64, and 128. We can observe that there is an obvious gap between 1-layer and 2-layer under all three width settings. However, if we deepen the number of layers, the gain is not remarkable. In light of these findings and taking into account the linear increase in computational overhead associated with additional layers, we apply two layers for both our small and medium configurations.

**The Length of Ray Feature Sequence.** We proceed to investigate the impact of feature sequence length across three distinct settings of LSTM width. The results of this exploration are systematically presented in Fig. 7 (b). A discernible trend emerges from the data, indicating that as the length of the feature sequence increases, there is a corresponding and gradual enhancement in image quality. To balance the computation afford and quality, we set 256 as the sequence length for Realistic Synthetic 360°.

## 5. Conclusion

In this paper, we proposed G-NeLF, a completely implicit NeLF method for novel view synthesis. G-NeLF does not rely on any physical modeling techniques such as volume density. Our method demonstrated that algorithms based on volume rendering, *e.g*., NeRF, aren't the sole avenues to address this challenge. Extensive experiments demonstrate that our method, even with a modest model size of 0.95 MB and utilizing only 100 views, still can successfully surpass the SoTA NeLF, R2L [40], training within approximately one hour. We hope that the comprehensive enhancements introduced in our method can contribute to the advancement of research in this area.

**Future Work:** It remains a promising research direction to extend our method to cater to dynamic scenes.

# References

[1] Benjamin Attal, Jia-Bin Huang, Michael Zollhöfer, Johannes Kopf, and Changil Kim. Learning neural light fields with ray-space embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19819–19829, 2022. 1, 2, 6, 7

[2] Haotian Bai, Yiqi Lin, Yize Chen, and Lin Wang. Dynamic plenoctree for adaptive sampling refinement in explicit nerf. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8785–8795, 2023. 3

[3] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 1, 2, 7

[4] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2023. 3

[5] Junli Cao, Huan Wang, Pavlo Chemerys, Vladislav Shakhrai, Ju Hu, Yun Fu, Denys Makoviichuk, Sergey Tulyakov, and Jian Ren. Real-time neural light field on mobile devices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8328–8337, 2023. 1, 2, 3, 4, 7

[6] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16123–16133, 2022. 3

[7] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022. 1, 2, 3, 6, 7, 8

[8] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16569–16578, 2023. 1, 3

[9] Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014. 8

[10] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 7, 2

[11] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. 2, 3

[12] Quankai Gao, Qiangeng Xu, Hao su, Ulrich Neumann, and Zexiang Xu. Strivec: Sparse tri-vector radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 3

[13] Quankai Gao, Qiangeng Xu, Hao Su, Ulrich Neumann, and Zexiang Xu. Strivec: Sparse tri-vector radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17569–17579, 2023. 1

[14] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021. 1, 3

[15] 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 - SIGGRAPH '96*, 1996. 2, 4

[16] Aarush Gupta, Junli Cao, Chaoyang Wang, Ju Hu, Sergey Tulyakov, Jian Ren, and László Jeni. Lightspeed: light and fast neural light fields on mobile devices. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 2, 4, 7

[17] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 2, 5, 8

[18] Wenbo Hu, Yuling Wang, Lin Ma, Bangbang Yang, Lin Gao, Xiao Liu, and Yuewen Ma. Tri-miprf: Tri-mip representation for efficient anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19774–19783, 2023. 1

[19] Lutao Jiang, Ruyi Ji, and Libo Zhang. Sdf-3dgan: A 3d object generative method based on implicit signed distance function. *arXiv preprint arXiv:2303.06821*, 2023. 3

[20] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 1, 4

[21] Andrej Karpathy, Justin Johnson, and Li Fei-Fei. Visualizing and understanding recurrent networks. *arXiv preprint arXiv:1506.02078*, 2015. 5, 8

[22] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques - SIGGRAPH '96*, 1996. 2, 4

[23] Lingzhi Li, Zhen Shen, Zhongshu Wang, Li Shen, and Liefeng Bo. Compressing volumetric radiance fields to 1 mb. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4222–4231, 2023. 1, 2, 3, 6, 7

[24] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14556–14565, 2021. 3

[25] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances*

*in Neural Information Processing Systems*, 33:15651–15663, 2020. 3

[26] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 1

[27] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 1

[28] 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 *European Conference on Computer Vision*, pages 405–421. Springer, 2020. 1, 2, 3, 4, 6, 7, 8

[29] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 1, 2, 3, 4, 5, 6, 7, 8

[30] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 1

[31] Lukas Radl, Andreas Kurz, and Markus Steinberger. Analyzing the internals of neural radiance fields. *arXiv preprint arXiv:2306.00696*, 2023. 1

[32] Vincent Sitzmann, Semon Rezchikov, Bill Freeman, Josh Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems*, 34: 19313–19325, 2021. 1, 2, 4

[33] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8269–8279, 2022. 1, 2, 4

[34] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 7, 2

[35] Towaki Takikawa, Alex Evans, Jonathan Tremblay, Thomas Müller, Morgan McGuire, Alec Jacobson, and Sanja Fidler. Variable bitrate neural fields. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022. 3

[36] Jiaxiang Tang, Xiaokang Chen, Jingbo Wang, and Gang Zeng. Compressible-composable nerf via rank-residual decomposition. *Advances in Neural Information Processing Systems*, 35:14798–14809, 2022. 3

[37] Jiaxiang Tang, Hang Zhou, Xiaokang Chen, Tianshu Hu, Errui Ding, Jingdong Wang, and Gang Zeng. Delicate textured mesh recovery from nerf via adaptive surface refinement. *arXiv preprint arXiv:2303.02091*, 2023. 3

[38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia

Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2, 5, 8

[39] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022. 5

[40] Huan Wang, Jian Ren, Zeng Huang, Kyle Olszewski, Menglei Chai, Yun Fu, and Sergey Tulyakov. R2l: Distilling neural radiance field to neural light field for efficient novel view synthesis. In *European Conference on Computer Vision*, pages 612–629. Springer, 2022. 1, 2, 4, 6, 7, 8

[41] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenoctrees for dynamic radiance field rendering in real-time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13524–13534, 2022. 3

[42] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. *arXiv preprint arXiv:2212.05231*, 2022. 3

[43] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf––: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 1

[44] Zhongshu Wang, Lingzhi Li, Zhen Shen, Li Shen, and Liefeng Bo. 4k-nerf: High fidelity neural radiance fields at ultra high resolutions. *arXiv preprint arXiv:2212.04701*, 2022.

[45] Wenqi Yang, Guanying Chen, Chaofeng Chen, Zhenfang Chen, and Kwan-Yee K Wong. Ps-nerf: Neural inverse rendering for multi-view photometric stereo. In *European Conference on Computer Vision*, pages 266–284. Springer, 2022. 1

[46] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 3

[47] Heng Yu, Joel Julin, Zoltan A Milacski, Koichiro Niinuma, and László A Jeni. Dylin: Making light field networks dynamic. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12397–12406, 2023. 1, 2, 4

[48] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, pages 11121–11128, 2023. 5, 8

[49] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 1

[50] Shuhong Zheng, Zhipeng Bao, Martial Hebert, and Yu-Xiong Wang. Multi-task view synthesis with neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21538–21549, 2023. 1

# G-NeLF: Memory- and Data-Efficient Hybrid Neural Light Field for Novel View Synthesis

## Supplementary Material

## 6. More comparison results

We show SSIM comparison results in the following tables. Due to the SSIM calculation bug in MobileR2L [5] and LightSpeed [16], we don't provide their SSIM. Please refer to the PSNR report for comparison.

| Method | AVG | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship | Size(MB) | Data Size |
|---|---|---|---|---|---|---|---|---|---|---|---|
| NeRF [28] | 0.964 | 0.976 | 0.945 | 0.968 | 0.983 | 0.973 | 0.967 | 0.974 | 0.926 | 4.56 | 100 |
| R2L [40] | 0.971 | 0.991 | 0.951 | 0.966 | 0.991 | 0.978 | 0.975 | 0.980 | 0.934 | 22.62 | 10,100 |
| Ours-S | 0.971 | 0.986 | 0.936 | 0.974 | 0.990 | 0.983 | 0.976 | 0.989 | 0.931 | 0.95 | 100 |
| Ours-M | 0.979 | 0.992 | 0.951 | 0.981 | 0.993 | 0.991 | 0.987 | 0.993 | 0.945 | 1.41 | 100 |
| Ours-L | 0.981 | 0.992 | 0.954 | 0.986 | 0.994 | 0.993 | 0.988 | 0.993 | 0.951 | 7.16 | 100 |

Table 9. SSIM quantitative comparison at $400 \times 400$ resolution on Realistic Synthetic 360° dataset. *Size* is model size. *Data Size* is the number of the training images. value means first, value means second, value means third.

| Method | Avg | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship | Size(MB) |
|---|---|---|---|---|---|---|---|---|---|---|
| VQRF [23] | 0.954 | 0.974 | 0.928 | 0.977 | 0.978 | 0.973 | 0.945 | 0.982 | 0.877 | 1.43 |
| NeRF [28] | 0.947 | 0.967 | 0.925 | 0.964 | 0.974 | 0.961 | 0.949 | 0.980 | 0.856 | 5.00 |
| Mip-NeRF [3] | 0.961 | 0.981 | 0.932 | 0.980 | 0.982 | 0.978 | 0.959 | 0.991 | 0.882 | 2.50 |
| TensoRF [7] | 0.963 | 0.985 | 0.937 | 0.982 | 0.982 | 0.983 | 0.952 | 0.988 | 0.895 | 71.80 |
| I-NGP [29] | 0.963 | 0.979 | 0.937 | 0.981 | 0.982 | 0.982 | 0.951 | 0.990 | 0.896 | 64.10 |
| Ours-M | 0.956 | 0.973 | 0.928 | 0.975 | 0.982 | 0.979 | 0.972 | 0.990 | 0.870 | **1.41** |
| Ours-L | 0.964 | 0.982 | 0.932 | 0.980 | 0.984 | 0.984 | 0.975 | 0.992 | 0.886 | 7.16 |

Table 10. SSIM quantitative comparison at $800 \times 800$ resolution on Realistic Synthetic 360° dataset. The VQRF is based on DVGO [34]. *Size* is model size. *Data Size* is the number of training images. **The first part** is NeRF's compression method. **The second part** is the SoTA NeRF methods. **The third part** is the SoTA NeLF methods.

| Method | Dense | | | | | | Sparse | | | Additional Images |
|---|---|---|---|---|---|---|---|---|---|---|
| | Avg | Horns(54) | Trex(48) | Room(36) | Fortress(36) | Flower(29) | Leaves(21) | Orchids(21) | Fern(17) | |
| NeRF [28] | 0.917 | 0.940 | 0.947 | 0.979 | 0.958 | 0.923 | 0.844 | 0.794 | 0.899 | 0 |
| Attal *et al.* [1] | 0.956 | 0.955 | 0.959 | 0.979 | 0.954 | 0.934 | 0.847 | 0.766 | 0.850 | 0 |
| R2L [40] | 0.950 | 0.934 | 0.948 | 0.980 | 0.959 | 0.929 | 0.860 | 0.800 | 0.895 | 10,000 |
| Ours | 0.961 | 0.961 | 0.961 | 0.979 | 0.963 | 0.939 | 0.858 | 0.716 | 0.873 | 0 |

Table 11. SSIM quantitative comparison at $504 \times 378$ resolution on Real Forward-Facing. *Additional Images* is the required number of extra synthesized images. The number of original training views is in the brackets.

| Method | Dense | | | | | | Sparse | | |
|---|---|---|---|---|---|---|---|---|---|
| | Avg | Horns(54) | Trex(48) | Room(36) | Fortress(36) | Flower(29) | Leaves(21) | Orchids(21) | Fern(17) |
| NeRF [28] | 0.873 | 0.828 | 0.880 | 0.948 | 0.881 | 0.827 | 0.690 | 0.641 | 0.792 |
| Plenoxels [10] | 0.886 | 0.857 | 0.890 | 0.937 | 0.885 | 0.862 | 0.760 | 0.687 | 0.832 |
| TensoRF [7] | 0.889 | 0.859 | 0.890 | 0.946 | 0.889 | 0.859 | 0.746 | 0.655 | 0.816 |
| Ours | 0.889 | 0.862 | 0.897 | 0.948 | 0.881 | 0.856 | 0.713 | 0.547 | 0.735 |

Table 12. SSIM quantitative comparison at $1008 \times 756$ resolution on Real Forward-Facing. *Additional Images* is the required number of extra synthesized images. The number of original training views is in the brackets.