

Harnessing Low-Frequency Neural Fields for Few-Shot View Synthesis

Liangchen Song¹ Zhong Li² Xuan Gong¹ Lele Chen² Zhang Chen² Yi Xu² Junsong Yuan¹
¹University at Buffalo ²OPPO

Abstract

Neural Radiance Fields (NeRF) have led to breakthroughs in the novel view synthesis problem. Positional Encoding (P.E.) is a critical factor that brings the impressive performance of NeRF, where low-dimensional coordinates are mapped to high-dimensional space to better recover scene details. However, blindly increasing the frequency of P.E. leads to overfitting when the reconstruction problem is highly underconstrained, e.g., few-shot images for training. We harness low-frequency neural fields to regularize high-frequency neural fields from overfitting to better address the problem of few-shot view synthesis. We propose reconstructing with a low-frequency only field and then finishing details with a high-frequency equipped field. Unlike most existing solutions that regularize the output space (i.e., rendered images), our regularization is conducted in the input space (i.e., signal frequency). We further propose a simple-yet-effective strategy for tuning the frequency to avoid overfitting few-shot inputs: enforcing consistency among the frequency domain of rendered 2D images. Thanks to the input space regularizing scheme, our method readily applies to inputs beyond spatial locations, such as the time dimension in dynamic scenes. Comparisons with state-of-the-art on both synthetic and natural datasets validate the effectiveness of our proposed solution for few-shot view synthesis. Code is available at <https://github.com/lsongx/halo>.

1. Introduction

Neural Radiance Field (NeRF) [29] and its extensions have shown promising results for novel view synthesis. In NeRF, coordinate-based multilayer perceptrons (MLPs), also referred to as neural fields, are adopted for continuously representing the geometry and appearance of the 3D scene. A critical factor that leads to the success of NeRF is Positional encoding (P.E.), which maps the low-dimensional coordinates to high-dimensional embeddings for representing the high-frequency details.

However, NeRF may generate unsatisfactory results when only a few posed images are available for training

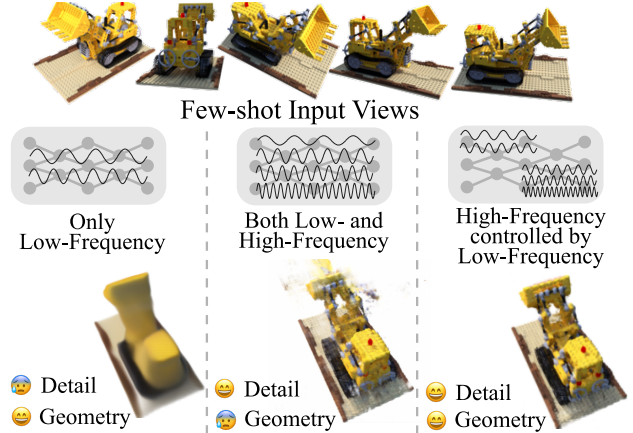


Figure 1. When solving the few-shot novel view synthesis problem with NeRF [29], the frequency of positional encoding (P.E.) has a significant impact on its scene representation performance. Adopting a low-frequency P.E. for input coordinates leads to smooth geometry with fewer details, while adopting P.E. with both low- and high-frequency leads to overfitting and erroneous geometry. We propose to harness low-frequency neural fields for regularizing the high-frequency neural fields.

[54]. Under the few-shot view synthesis setting, the original NeRF overfits the input views and converges to a degenerate solution [15]. When the radiance field is underconstrained, the high-dimensional P.E. in NeRF tends to fill the space with high-frequency geometry, resulting in poor generalization to unseen views. On the other hand, if we learn the neural field with low-frequency only P.E., fine details will be lost in the rendered views. Much significant progress has been made in adjusting the overall frequency of P.E. for a scene [45, 26, 14, 24, 2], however, in this paper, we argue that the potential behind low-frequency only neural fields can be further exploited for the few-shot setting.

We are motivated by the difference between low-frequency only and high-frequency added neural fields on interpolation and extrapolation of novel input coordinates. Low-frequency only neural fields generate smooth interpolation but fail to extrapolate high-frequency periodic signals. In contrast, neural fields with sufficient frequency can adapt to periodic patterns but bring high-frequency interpolation artifacts (Fig. 2). An idea then arises that we can do

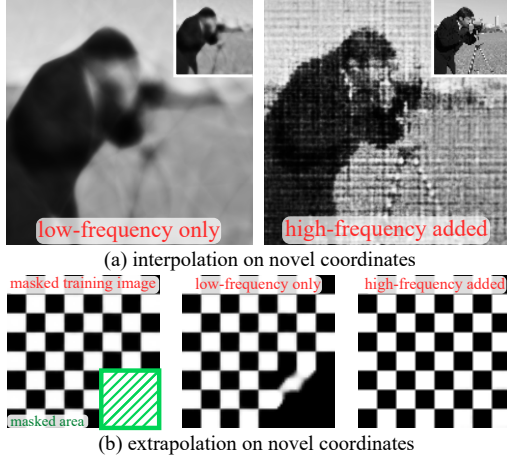


Figure 2. A 2D toy example for demonstrating different interpolation and extrapolation results between a low-frequency only neural field and a high-frequency added neural field. The neural fields take as input the 2D coordinates of pixels and predict the intensities. (a) Given a 64^2 image, the neural fields interpolate it to a 256^2 image. The top right corner demonstrates the prediction on the training image after 10000 iterations. The low-frequency only neural field interpolates smooth values while the high-frequency equipped neural field generates obvious structured artifacts. (b) Given a checkerboard image with a corner masked, the neural fields extrapolate the masked area. The low-frequency only neural field cannot extrapolate the pattern to unobserved regions while the high-frequency added neural field succeeds.

interpolation with low-frequency only neural fields and extrapolation with high-frequency added neural fields. However, turning the idea into a practical algorithm is challenging since it involves classifying a prediction on novel coordinates as an interpolation or an extrapolation.

For the problem of few-shot synthesis, we treat the rough geometry prediction as the low-frequency interpolation only problem and the fine detail prediction as the high-frequency extrapolation needed problem. Our general idea is demonstrated in Fig. 1. Since querying the geometry information from point-based fields is expensive, we design a framework consisting of a low-frequency only ray-based field for modeling rough geometry and a high-frequency added point-based field for modeling fine details. The ray-based field is adapted from recent proposed light field networks (LFN) [43], which has been proved far more efficient than point-based fields. During inference, the low-frequency only ray-based field will predict a depth value for each ray, instead of colors as in LFN. In this way, not only can the ray-based be trained easier since it only needs to fit a low-frequency smooth signal, but also the point-based radiance field can be efficiently regularized.

Furthermore, we propose to leverage the well-established frequency analysis of 2D image signals for determining the frequency of modeling rough geometry. Pre-

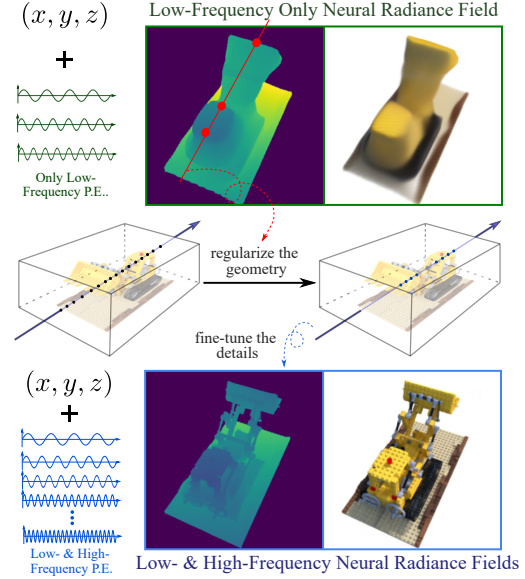


Figure 3. Controlling the geometry with a low-frequency neural field. The empty space of the two fields is enforced to be similar.

vious methods determine the frequency of neural fields by human or ground-truth-involved approaches [55], which is not applicable in our few-shot view synthesis setting. We assume that slightly changing the viewpoint should not significantly impact the frequency domain of rendered 2D images if the scene is free of high-frequency floating artifacts. Based on the assumption, we set a threshold on the change in the frequency domain and then tune the frequency of P.E. from high to low until the rendered images satisfy the predefined threshold. The proposed criterion is simple but effective and does not require extra data.

For validating our proposed method, we first follow the 8 views setting and the 14 one side view setting used by Diet-NeRF [15]. On this 360-degree rendering task, we train a low-frequency radiance field with the given eight views, and then it is used for supervising a ray-based field. Next, we consider a setting of 4 corner input views with challenging baselines on real data, which is practically appealing due to the complex geometry. In summary, our contributions are as follows:

- For few-shot novel view synthesis, we propose to harness the low-frequency inputs for regularizing the geometry, and then high-frequency inputs are added for details. Our regularization is conducted in the input space, enabling few-shot view synthesis on both static and dynamic scenes.
- Benefiting from the different interpolation and extrapolation results between low-frequency only and high-frequency added neural fields, our framework can extrapolate periodic patterns in unobserved regions.

- We propose a simple-yet-effective approach for tuning the frequency of P.E. under the few-shot view synthesis setting. The approach does not require extra data or labels, and the only hyperparameter (the threshold for the difference in rendered images’ frequency domain) is well-generalizable across different scenes.

2. Related Work

Neural fields and the frequency of inputs. Neural fields are also known as implicit neural representations or coordinate-based representations. Implicit neural representations are initially used for representing the geometry [27, 39, 33, 5, 6, 16], but few works are concerned with high-frequency modeling details with a neural network. A recent milestone work in the field of novel view synthesis is NeRF [29], in which a 5D field represents the scene. NeRF’s impressive performance is largely due to the positional encoding module, which maps the input 5D coordinate to a higher dimensional space. The high-frequency mapping scheme is theoretically studied by FourierFeat [45] and SIREN [42], in which the interpolation and extrapolation behavior with different input frequencies are explored as well. For stabilizing training, the smooth prediction of a low-frequency prediction is leveraged in [34, 23]. Coarse-to-fine frequency hierarchy is further adopted by Acorn [26], BACON [24], SAPE [14] and [53]. The previous methods also demonstrate the benefits of low-frequency neural fields, but harnessing low-frequency neural fields has yet to be well exploited under the setting of few-shot view synthesis.

NeRF for few-shot view synthesis. A limitation of NeRF is that a certain number of input views are required for constructing a good radiance field. One reason is that the neural field is optimized on each scene separately and many excellent works are introducing extra supervisory signals for few-shot input views. PixelNeRF [54] proposes to condition the radiance field on the semantics of the input views. Learning a latent code for each scene has also been adopted by [40, 48, 13]. GRAF [40] uses a discriminator on the rendered novel view patches. In [44], meta-learning is employed for learning a prior on the radiance field. Similarly, meta-learning has also been adopted for SDF [41] and light field [43, 10]. Semantic priors learned from large-scale image or language datasets like CLIP are adopted as supervisory signals in DietNeRF [15]. Depth or sparse point clouds reconstructed by classical methods are adopted in DSNeRF [8], MVSNeRF [4], NerfingMVS [50] and [38]. InfoNeRF [19] uses the entropy constraint of the density in each ray to minimize potential reconstruction inconsistency. RegNeRF [32] proposes regularizing the geometry and appearance of patches rendered from unobserved viewpoints. Different from most of the existing few-shot NeRF methods that regularize the model *in the output space* (i.e., the ren-

dered images), our regularization is conducted *in the input space*. Concurrent work FreeNeRF [52] is also inspired by the phenomenon of overfitting with high-frequency inputs. They propose to associate the frequency with visible ratio, while in our work we first train with low-frequency only and then adding high-frequency. The simplicity of our method enables readily apply the method to dynamic scenes that have not been studied in great detail.

Light field rendering with few-shot inputs. Light field rendering is a widely studied and applied technology [20, 12], and the angular-spatial resolution tradeoff is a long-standing problem. In [18], the authors propose a learning-based approach for rendering novel views with four corner sub-aperture views from the light fields captured by the Lytro Illum camera. The concept of multi-plane images (MPIs) [57] is proposed for decomposing the observation views and learning view extrapolation. Predicting the depth of images for rendering novel views is also studied in [7, 47, 51] with few inputs or only one input view. While there is much progress on the few-shot inputs setting, there is much room for improving the rendering quality without using a different dataset or pre-trained network.

3. Preliminaries

In NeRF, a neural field F_{Θ} parameterized by an MLP that takes 3D point locations $\mathbf{p} = (x, y, z)$ and viewing direction \mathbf{d} as input. The field is a mapping from each 3D location to its attributions $(\mathbf{c}, \sigma) = F_{\Theta}(\mathbf{p}, \mathbf{d})$, where $\mathbf{c} = (r, g, b)$ is the color and σ is volume density. For rendering, each pixel color is acquired from points along the ray through the field F_{Θ} and differentiable volume rendering. Each ray emitted from the camera center \mathbf{o} with direction \mathbf{d} , points along the ray are $\mathbf{p} = \mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ and the expected color $C(\mathbf{r})$ with near and far bounds t_n, t_f is $C(\mathbf{r}) = \int_{t_n}^{t_f} e^{-\int_{t_n}^t \sigma(\mathbf{r}(s))ds} \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt$. A set of points are sampled along the ray for numerically estimating the integration. During training, the interval $[t_n, t_f]$ is partitioned into K bins and for i th bin we uniformly sample a t with $t_i \sim U(t_n + \frac{i-1}{N}(t_f - t_n), t_n + \frac{i}{N}(t_f - t_n))$. Then the color of the ray is computed with

$$\hat{C}(\mathbf{r}) = \sum_{k=1}^K e^{-\sum_{k'=1}^{k-1} \sigma_{k'} \delta_{k'}} (1 - e^{-\sigma_k \delta_k}) \mathbf{c}_k, \quad (1)$$

where $\sigma_k = \sigma(\mathbf{r}(t_k))$, $\mathbf{c}_k = \mathbf{c}(\mathbf{r}(t_k))$ and $\delta_k = t_{k+1} - t_k$ is the distance between the two samples.

The points $\mathbf{p} \in \mathbb{R}^3$ are low-dimensional and cannot represent the high-frequency details of the scene. In NeRF, the coordinate inputs to a higher dimension with positional encoding

$$\gamma(\mathbf{p}) = (\sin(\mathbf{p}), \cos(\mathbf{p}), \dots, \sin(2^{L-1}\mathbf{p}), \cos(2^{L-1}\mathbf{p})), \quad (2)$$

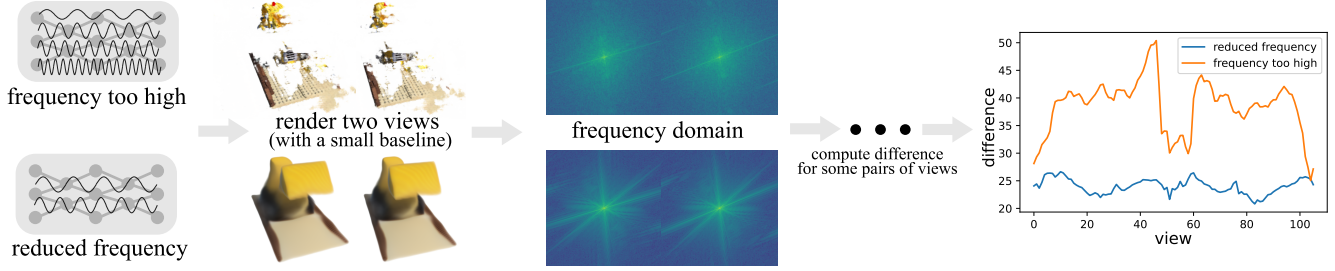


Figure 4. The proposed criteria for tuning the frequency of NeRF. A set of pairs of rendered views with a small baseline are considered. The difference of the pairs in the frequency domain is adopted as an indicator.

where L is a hyperparameter. Positional encoding is critical for NeRF, but a high L will lead to poor interpolation results, which has been theoretically studied by Tancik *et al.* [45]. Two fields, coarse and fine fields, are designed and trained jointly in NeRF. Low- and high-frequency P.E. are used for both two fields. The fine field takes samples near the visible contents predicted by the coarse field. Finally, the two fields are optimized with the reconstruction loss,

$$L_{\text{rec}} = \sum_{\mathbf{r}} \|\hat{C}(\mathbf{r}) - C_{\text{gt}}(\mathbf{r})\|_2^2, \quad (3)$$

where C_{gt} is the ground truth RGB colors for ray \mathbf{r} .

4. Proposed Method

The core idea of our method is to control the high-frequency components with a low-frequency only field. We first present a criterion for determining the frequency of Lo-NeRF, then introduce a framework for efficiently regularizing high-frequency. For simplicity, we denote the low-frequency only P.E. as *Lo* and high-frequency equipped P.E. as *Hi*. Lo-P.E. and Hi-P.E. equipped NeRF are denoted as Lo-NeRF and Hi-NeRF.

4.1. Frequency tuning

Since we are concerned with the few-shot view synthesis setting, we design a criterion based on the rendered images' frequency domain and do not require extra data or labeling. An overview of the criteria is demonstrated in Fig. 4. Let N paired images with a small baseline, denoted as $\{I_i^a, I_i^b\}$, are randomly rendered. Then the averaged difference on the pairs in the frequency domain is $\sigma = \frac{1}{N} \sum_i \|\mathcal{F}(I_i^a) - \mathcal{F}(I_i^b)\|$, where $\mathcal{F}()$ is Fourier transformation. Note that a mask is applied on the frequency domain to remove large values in the frequency domain for robustness (and detailed implementation can be found in supplementary). We reduce the frequency of NeRF until σ reaches a predefined threshold (25 in our experiment). The frequency is then used for defining Lo-NeRF.

4.2. Efficient low-frequency based regularizer

We demonstrate in Fig. 1 that Lo-NeRF produces a smooth geometry of the scene. A straightforward way of leveraging the smooth geometry from Lo-NeRF is separating the training into two stages: after training the Lo-NeRF, we enforce that the Hi-NeRF shares a similar rough geometry as the Lo-NeRF. The approach is illustrated in Fig. 3.

Low-frequency ray-based field. We propose to sample random rays in the space and then enforce the Hi-NeRF having consistent geometry for these random rays. However, using a Lo-NeRF for regularizing can be time-consuming since querying the depth of a ray is rather expensive. Thus, we adopt a ray-based field for regularizing the Hi-NeRF in an online manner: The computational cost can be reduced by directly predicting the depth value for each ray, as in DOnERF [30]. Ray-based predictions may suffer from blurry edges in DOnERF. However, blurry predictions are acceptable in our case since only rough geometry is needed for regularization, and Hi-NeRF can fine-tune the details.

For training, the depth of a ray is first computed by sampling points and accumulates density outputs from the low-frequency point-based field. Then a ray-based field is supervised with the depth. All the ray origins \mathbf{o} are the intersection between the ray at a predefined sphere and view direction \mathbf{d} are unit vectors. So though the inputs of the ray-based field (\mathbf{o}, \mathbf{d}) are six-dimensional, the degree of freedom is four since they are points from two surfaces. Denote the ray-based field as F_{ray} and the point-based field (*i.e.*, NeRF) as F_{point} , then the loss for training F_{ray} is

$$L_{\text{ray}} = \sum_{\mathbf{r}=(\mathbf{o}, \mathbf{d})} \|F_{\text{ray}}(\mathbf{o}, \mathbf{d}) - D(\mathbf{r}; F_{\text{point}})\|_2^2, \quad (4)$$

where $D(\mathbf{r}; F_{\text{point}})$ is the depth computed with F_{point} .

Empty space loss. To regularize the empty space (*i.e.*, where depth cannot be inferred from Lo-NeRF), we add an extra loss on the accumulated occupancy for a random ray \mathbf{r} , which is defined as $\text{acc}(\mathbf{r}) = \sum_{k=1}^K e^{-\sum_{k'=1}^{k-1} \sigma_{k'} \delta_{k'}} (1 -$

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF [29]	14.934	0.687	0.318
NV [25]	17.859	0.741	0.245
Simplified NeRF [15]	20.092	0.822	0.179
DietNeRF [15]	23.147	0.866	0.109
DietNeRF ft [15]	23.591	0.874	0.097
HALO (<i>Ours</i>)	<u>23.269</u>	<u>0.863</u>	0.152
HALO+DietNeRF	23.687	0.879	0.090
NeRF, 100 views	31.153	0.954	0.046

Table 1. Eight training views are randomly sampled for each scene on Realistic Synthetic scenes. Metrics averaged across the 8 scenes are reported.

$e^{-\sigma_k \delta_k}$). Denote the accumulated occupancy for the Lo-NeRF and Hi-NeRF as acc_{Lo} and acc_{Hi} respectively, then the empty space regularization loss is

$$L_{\text{empty}} = \sum_{\{\mathbf{r}: \text{acc}_{\text{Lo}}(\mathbf{r}) < \tau\}} \text{acc}_{\text{Hi}}(\mathbf{r}), \quad (5)$$

where τ is a threshold for determining if the space is empty, and we set $\tau = 0.01$. Since accumulated occupancy values are all larger than 0, simply summing all the values up can effectively regularize the occupancy of all points on the ray to be 0.

4.3. Overall framework

As described above, our overall framework consists of 3 stages: Lo-NeRF training, ray-based field training, and Hi-NeRF training. Lo-NeRF training is the same as training a typical NeRF except that the frequency of P.E. (*i.e.*, L in Eq. (2)) is adjusted according to the frequency domain of rendered images. During the second stage, *i.e.*, ray-based field training, we freeze the Lo-NeRF, and then random initialize a ray-based field. The ray-based field is then optimized with L_{ray} from Eq. (4). For the final stage, all previous networks are frozen, and we random initialize another NeRF with the default frequency in NeRF ($L = 10$). The Hi-NeRF to be optimized will sample points according to the depth predicted by the ray-based field, following the same sampling strategy as in DNeRF. The training loss for the Hi-NeRF then becomes

$$L = L_{\text{rec}} + \lambda L_{\text{empty}},$$

where λ is a balancing parameter. L_{rec} is computed by Eq. (3) and L_{empty} is computed by Eq. (5).

5. Experiments

We validate our method in two aspects. First, we follow the few-shot settings in DietNeRF [15]: 360° rendering from only 8 views and rendering on views not observed during training. Next, we present rendering results on a two-plane light field dataset with 4 corner views. Second, we

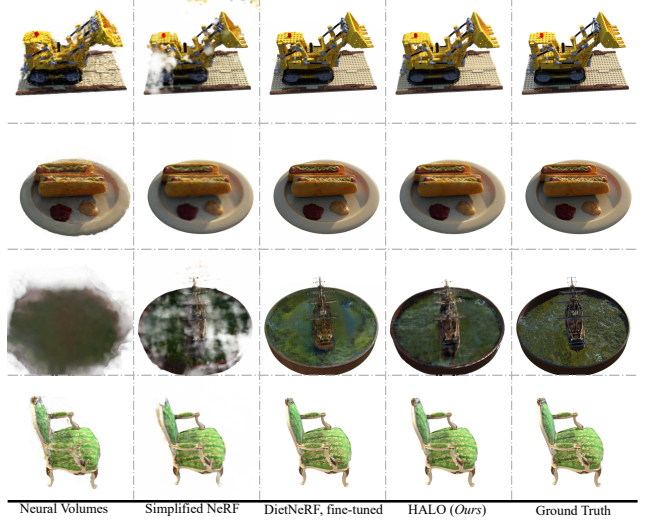


Figure 5. Comparison of novel views rendering results under the random 8 views setting.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF [29]	19.662	0.799	0.202
Simplified NeRF [15]	21.553	0.818	0.160
DietNeRF [15]	20.753	0.810	0.157
DietNeRF ft [15]	<u>22.211</u>	<u>0.824</u>	0.143
HALO (<i>Ours</i>)	22.581	0.827	<u>0.150</u>
HALO+DietNeRF	22.862	0.833	0.140
NeRF, 100 views	31.618	0.965	0.033

Table 2. A quantitative comparison under the one side views setting for the “Lego” scene. Note that DietNeRF requires a trained CLIP model [37] for supervision, while we do not.

test our method on dynamic scenes from D-NeRF [36] with every 8 image from the original training set. Our method is denoted as **HALO** (short for **H**arnessing **L**ow-frequency). Detailed implementation is described in the supplementary.

5.1. Datasets and evaluation metrics.

For 360° rendering tasks, we use the Realistic Synthetic scenes [29]. There are 8 objects and images rendered from the objects are split into the train, validation, and test sets. We report results on the test set. Stanford Light Field Archive (StanfordLF) [1] is used to evaluate the two-plane light field data. There are 12 scenes, and each scene has 17×17 images. For each scene, we select the images with index (4,4) (4,12) (12,4) (12,12) as the four corner images. Rendering results on images indexed by (8,6) (8,10) (6,8) (10,8) are compared to ground truth for evaluation. The LLFF [28, 29] dataset is also used for evaluating the performance of non-structured light field data. We manually select 4 views from 4 challenging scenes in LLFF as the training set, and an image inside the 4 views is used for evaluation. We use the DNeRF dataset [36] and every eighth

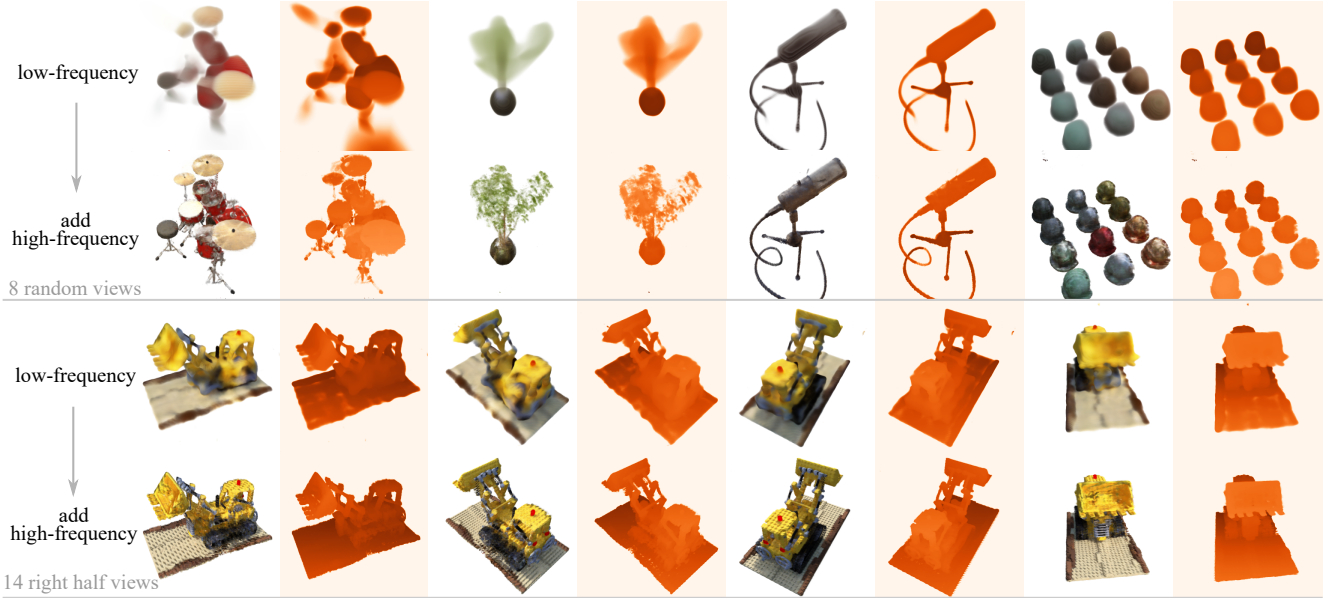


Figure 6. Novel views synthesized from the low-frequency neural field and the corresponding controlled high-frequency neural field.

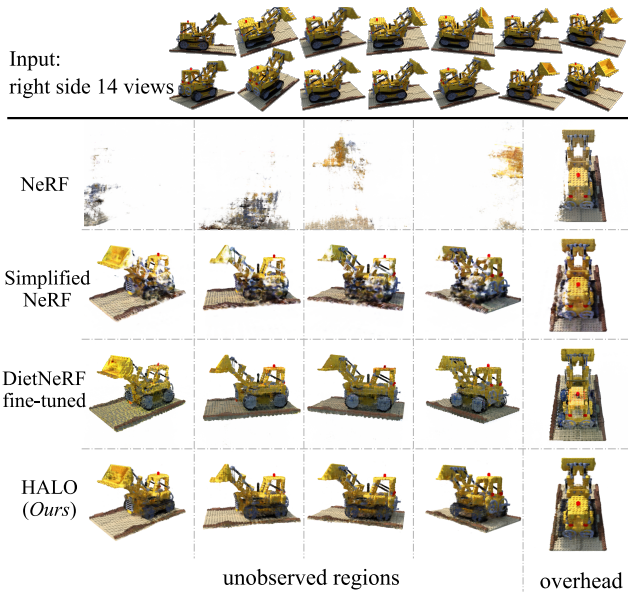


Figure 7. Novel views extrapolation with right side 14 views as inputs. The testing views are from the other side, which is not covered by the training views.

training image as the few-shot inputs for dynamic scenes. Different from static scenes, few-shot inputs on DNeRF indicate being sparse both spatial and temporal.

For evaluation, three metrics are used for evaluating the rendering results: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) [49] and Learned Perceptual Image Patch Similarity (LPIPS) [56] with VGG backbone.

5.2. 360° scenes

We validate our method with the 360° rendering settings proposed by DietNeRF [15] in this section.

Random views as inputs To test the few-shot performance, 8 views are randomly selected by DietNeRF. We use the same random views for testing (detailed image names included in the supplementary). Overall performance on all scenes are reported in Tab. 1. Our results achieve comparable results to DietNeRF without using extra semantic supervision signals. Fig. 5 visualizes the reconstruction results of the methods. An interesting point can be observed by comparing our method and DietNeRF on the ship object. DietNeRF-ft and our method generate decent results, but the brightness of our reconstruction is more consistent with the ground truth. This demonstrates that a potential drawback of introducing semantic supervision affects the lighting consistency, while our HALO can fill in the gap. Further, our method and DietNeRF are orthogonal techniques and can potentially be combined. We report in the table a naive combination that first training with our method and then adding DietNeRF based fine tuning. The improvements demonstrate the potential of combining our method and other data-driven prior based methods.

One side views as inputs In DietNeRF, the authors demonstrate the extrapolation ability enabled by introducing semantic supervision through reconstructing unobserved regions. We follow the setting of using 14 input views from the right side, and results are reported in Tab. 2 and Fig. 7. It is intriguing to observe that a neural field it-

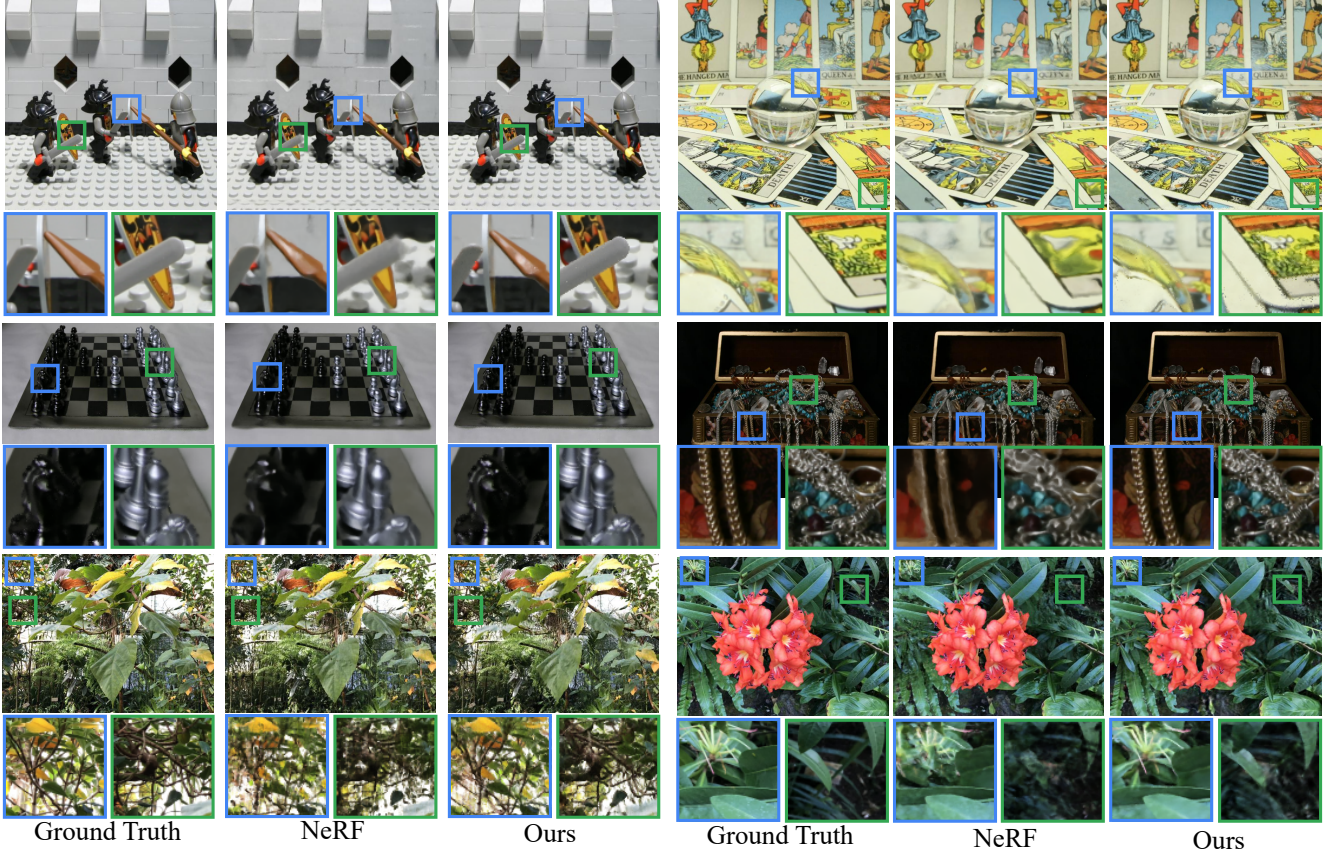


Figure 8. Qualitative results on StanfordLF (top & middle rows) and LLFF (bottom row). The training views are four corner views, and the testing views are the center view. Despite being few-shot (4 views), the baseline is small, so NeRF can render decent images. The results demonstrate that harnessing low-frequency can improve details on few-shot inputs with small baselines.

self can well extrapolate the complex periodic signals since no extra supervisory signals are adopted in our method. Moreover, our method outperforms DietNeRF-ft in terms of PSNR and SSIM. By comparing the results in Fig. 7, we observe similar effects as the ship scene in Fig. 5: the brightness of images reconstructed by DietNeRF is not consistent with the ground truth. This explains the reason behind the higher LPIPS score achieved by DietNeRF as the perceptual metric cares less about colors [17].

We also demonstrate the outputs from low-frequency and high-frequency neural fields in Fig. 6. We can observe from the images that the extrapolation behavior between ours and DietNeRF is quite different: DietNeRF generates structurally consistent new contents, while ours extrapolate periodic contents of both structure and texture with high fidelity. Our method’s artifacts are mainly on edges next to the wheels, and we attribute them to the complex structural signal with poor periodic patterns.

5.3. Forward-facing scenes

We evaluate our framework on real-world forward-facing images in this section. Four corner views are selected

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
<i>On Stanford Lightfield Archive [1]</i>			
NeRF [29]	30.487	0.819	0.248
HALO (Ours)	31.283	0.897	0.234
<i>On LLFF [29]</i>			
NeRF [29]	19.907	0.634	0.340
HALO (Ours)	20.578	0.627	0.322

Table 3. Quantitative comparison on forward-facing scenes.

as the training view for each scene. Note that the camera intrinsics and extrinsics are not provided in StanfordLF, so to test the performance of NeRF on the dataset StanfordLF, we use an EPI-based parameterization in NeRF. A more detailed definition of the coordinates is included in the supplementary. Besides the two-plane light field dataset, we also compare our method with NeRF on the LLFF dataset [28]. In Tab. 3, we compare our method with NeRF. Fig. 8 visualizes the rendered images from our method and NeRF. The top two rows of images are from the StanfordLF dataset, and the third row is from the LLFF dataset. The images in this section are with small baselines. Therefore the vanilla

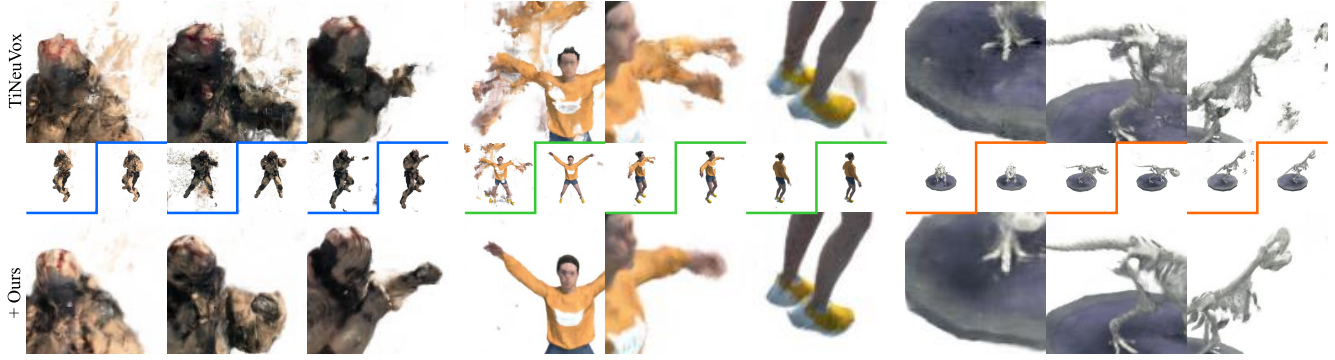


Figure 9. Qualitative comparisons on dynamic scenes. TiNeuVox [9] uses both low- and high-frequency by default, thus overfitting few-shot inputs. After adopting our method, the images have finer details and fewer floating outliers.

Method	Time (GPU/h) ↓	PSNR↑	SSIM↑	LPIPS↓
TiNeuVox-S [9]	0.17	21.392	0.864	0.182
DietNeRF [15]	4.25	22.487	0.873	0.177
DietNeRF-ft [15]	5.17	22.815	0.878	0.171
RegNeRF [31]	3.33	23.995	0.889	0.143
HALO (Ours)	0.67	25.687	0.914	0.116

Table 4. Quantitative comparisons on dynamic scenes. Our method can be easily extended to extra input dimension (*i.e.*, time t), while other methods are applied frame-by-frame. Time indicates the training time on one scene. Our method does not require querying a pre-trained prior network, thus keeping efficient.

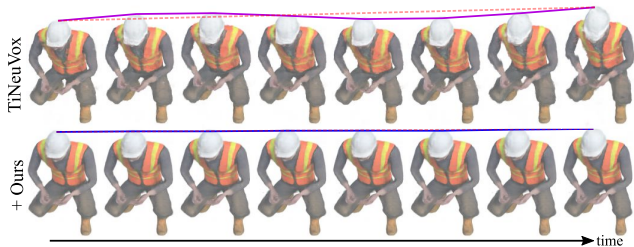


Figure 10. Temporal flickering motion is observed on TiNeuVox under few-shot dynamic settings since the default frequency of time t is high. Solid lines demonstrate the trajectory and dashed lines demonstrate the start-to-end linear trajectory.

NeRF can generate decent novel views despite the access to only 4 corner views. We can observe from Tab. 3 and Fig. 8 that NeRF will not overfit the training views and no floating outliers are presented with these inputs, harnessing low-frequency can still be beneficial for rendering realistic novel views.

5.4. Dynamic scenes

Our method uses low-frequency inputs as a regularization; thus, extra dimension in the input space can be readily implemented. We consider an extra time dimension and test our method on dynamic scenes. For dynamic scenes, not only are multiview observations sparse, sampling on the time axis becomes sparse as well for few-shot challenges. We adopt the recent efficient reconstruction method

TiNeuVox [9] as the baseline, which can be trained within minutes. Like in vanilla NeRF, low-frequency reconstruction can be achieved by reducing the frequency in positional encoding. We skip training a low-frequency ray-based field in this experiment since directly querying from a low-frequency TiNeuVox is cheap. When tuning the low-frequency input, we use a similar strategy as on static scenes, but with a modification that renders new images with a fixed timestamp.

In Tab. 4, we present the quantitative results on the dynamic setting. TiNeuVox-S is adopted as the base model for all the experiments. TiNeuVox uses both low- and high-frequency for reconstruction by default, while low-frequency inputs regularize ours. We can observe that our method consistently improves performance when observations are sparse. We visualize rendered images in Fig. 9, from which we can observe clear visual improvements after adopting our regularization method. Our method does not require data-driven priors to regularize, therefore keeping the efficiency of volume based NeRF models. Furthermore, we observe that the high-frequency of the time axis introduces temporal flickering motion when observations of dynamic sequence are sparse. An illustration is demonstrated in Fig. 10, and more can be found in our video.

6. Conclusion

We propose to harness the low-frequency NeRF and leverage it to regularize the high-frequency NeRF so that it will not overfit under the few-shot setting. The regularization is conducted in the input space so our method can be readily applied to static and dynamic scenes. Also benefiting from the different interpolation and extrapolation properties of low- and high-frequency, our method can extrapolate periodic contents and render realistic novel views on unobserved areas. Furthermore, we design a simple-yet-effective criterion for determining a NeRF’s frequency to avoid overfitting. Our experimental results demonstrate the effectiveness of our proposed solution’s effectiveness and the potential of harnessing low-frequency neural fields.

References

- [1] Andrew Adams. The (new) stanford light field archive. <http://lightfield.stanford.edu/lfs.html>, 2008. 5, 7
- [2] Nuri Benbarka, Timon Höfer, Hamd ul-Moqet Riaz, and Andreas Zell. Seeing implicit neural representations as fourier series. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2041–2050, January 2022. 1
- [3] Robert C Bolles, H Harlyn Baker, and David H Marimont. Epipolar-plane image analysis: An approach to determining structure from motion. *International journal of computer vision*, 1(1):7–55, 1987. 12
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, October 2021. 3
- [5] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. 3
- [6] Julian Chibane, Gerard Pons-Moll, et al. Neural unsigned distance fields for implicit function learning. *Advances in Neural Information Processing Systems*, 33, 2020. 3
- [7] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7781–7790, 2019. 3
- [8] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. *arXiv preprint arXiv:2107.02791*, 2021. 3
- [9] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. *arxiv:2205.15285*, 2022. 8, 18
- [10] Brandon Yushan Feng and Amitabh Varshney. Signet: Efficient neural representation for light fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14224–14233, October 2021. 3
- [11] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. 12
- [12] Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. The lumigraph. In John Fujii, editor, *SIGGRAPH*, pages 43–54. ACM, 1996. 3, 12
- [13] Philipp Henzler, Jeremy Reizenstein, Patrick Labatut, Roman Shapovalov, Tobias Ritschel, Andrea Vedaldi, and David Novotny. Unsupervised learning of 3d object categories from videos in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4700–4709, 2021. 3
- [14] Amir Hertz, Or Perel, Raja Giryes, Olga Sorkine-hornung, and Daniel Cohen-or. SAPE: Spatially-adaptive progressive encoding for neural optimization. In *Advances in Neural Information Processing Systems*, 2021. 1, 3
- [15] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, October 2021. 1, 2, 3, 5, 6, 8
- [16] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6001–6010, 2020. 3
- [17] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 7
- [18] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia 2016)*, 35(6), 2016. 3
- [19] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022. 3
- [20] Marc Levoy and Pat Hanrahan. Light field rendering. In John Fujii, editor, *SIGGRAPH*, pages 31–42. ACM, 1996. 3, 12
- [21] Tianye Li, Mira Slavcheva, Michael Zollhofer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, and Zhaoyang Lv. Neural 3d video synthesis. *arXiv preprint arXiv:2103.02597*, 2021. 12
- [22] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 12
- [23] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. 3
- [24] David B Lindell, Dave Van Veen, Jeong Joon Park, and Gordon Wetzstein. Bacon: Band-limited coordinate networks for multiscale scene representation. *arXiv preprint arXiv:2112.04645*, 2021. 1, 3
- [25] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Transactions on Graphics*, 38(4):1–14, 2019. 5
- [26] 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 Trans. Graph.*, 40(4):58:1–58:13, 2021. 1, 3
- [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. 3
- [28] 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 (Proceedings of SIGGRAPH)*, 38(4), 2019. [5](#), [7](#), [14](#)
- [29] 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](#), [3](#), [5](#), [7](#)
- [30] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarthy R. Alla Chaitanya, Anton S. Kaplanyan, and Markus Steinberger. DONeRF: Towards Real-Time Rendering of Compact Neural Radiance Fields using Depth Oracle Networks. *Computer Graphics Forum*, 40(4), 2021. [4](#)
- [31] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. *arXiv preprint arXiv:2112.00724*, 2021. [8](#)
- [32] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. [3](#)
- [33] 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. [3](#)
- [34] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, October 2021. [3](#), [12](#)
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019. [12](#)
- [36] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. [5](#), [12](#), [14](#), [18](#)
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. [5](#)
- [38] Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. [3](#)
- [39] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2304–2314, 2019. [3](#)
- [40] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33, 2020. [3](#)
- [41] Vincent Sitzmann, Eric R. Chan, Richard Tucker, Noah Snively, and Gordon Wetzstein. Metasdf: Meta-learning signed distance functions. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing*, 2020. [3](#)
- [42] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33, 2020. [3](#)
- [43] Vincent Sitzmann, Semon Rezhikov, William T. Freeman, Joshua B. Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. In *Advances in Neural Information Processing Systems*, 2021. [2](#), [3](#)
- [44] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2846–2855, 2021. [3](#)
- [45] 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 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems*, 2020. [1](#), [3](#), [4](#), [14](#)
- [46] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, October 2021. [12](#)
- [47] Richard Tucker and Noah Snively. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 551–560, 2020. [3](#)
- [48] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snively, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. [3](#)
- [49] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [6](#)

- [50] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *ICCV*, 2021. 3
- [51] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7467–7477, 2020. 3
- [52] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. 2023. 3
- [53] Wang Yifan, Lukas Rahmann, and Olga Sorkine-hornung. Geometry-consistent neural shape representation with implicit displacement fields. In *International Conference on Learning Representations*, 2022. 3
- [54] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 1, 3
- [55] Huangjie Yu, Anpei Chen, Xin Chen, Lan Xu, Ziyu Shao, and Jingyi Yu. Anisotropic fourier features for neural image-based rendering and relighting. In *AAAI Conference on Artificial Intelligence*, pages 3152–3160. AAAI Press, 2022. 2
- [56] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6
- [57] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM Transactions on Graphics*, 37(4):1–12, 2018. 3

Appendix

A. Extend to light field data

To further validate the proposed framework, we consider the angular-spatial resolution tradeoff on two-plane parameterized light field data. Though we can estimate camera intrinsics and extrinsics from images, directly processing light rays is a more general approach and sometimes favorable. We design a parameterization for the 3D points in the space, which enables controlling the geometry with light ray inputs.

EPI-based parameterization. The two-plane parameterization light field [20, 12] is a well-established method for representing rays in the space. Each ray is parameterized by the intersection points on the camera plane uv and the image plane st . Then, a point in space will be a line on the us or st slice of the light field (Fig. 11), which is also known as Epipolar Plane Images (EPIs) [3]. We denote the slope of the line as $\tan\theta$, then a 3D point can be determined by (u, v, s, t, θ) . Consequently, points on the ray (u, v, s, t) can be represented by $\{(u, v, s, t, \theta_i)\}_{i=0}^N$, where N is the number of sampled points. To get the color of the ray, we take the integral over θ .

The parameterization (u, v, s, t, θ) for 3D points is redundant and can be simplified. In Fig. 12, we demonstrate how v and t are correlated: For the same 3D point, two rays observing the point will satisfy $\frac{\Delta v}{\Delta t} = \arctan \theta = \text{constant}$ as the two planes are fixed. Based on this fact, we propose to align all (u, v, s, t, θ) to a fixed uv for representing the spatial location of a 3D point. Let the u^* and v^* be the fixed value, then $s' = s + \frac{u-u^*}{\arctan \theta}$ and $t' = t + \frac{v-v^*}{\arctan \theta}$. In this way, (u, v, s, t, θ) and $(u^*, v^*, s', t', \theta)$ represent the same 3D point. After aligning all 3D points to the same u^*v^* , each point is now parameterized by a 3D vector (s', t', θ) . The inputs of radiance field are then switched to (s', t', θ) for 3D points and (u, v) for viewing direction.

Joint training The two-stage training scheme limits the framework’s extensibility since the low-frequency cannot be further optimized once finished. For example, a motion field which is widely used for modeling dynamic scenes (e.g., [11, 36, 22, 46, 21, 34]) cannot be directly adopted into the two-stage framework.

An overview of our proposed framework for harnessing low-frequency neural fields is illustrated in Fig. 13. Inputs from 4D light field data are used for demonstration. For rendering with general 3D world coordinates (i.e., with known intrinsics and extrinsics), there are two small differences: First, the inputs for points are now (x, y, z) and the align procedure is not needed; Second, the target for L_{consist} is from another Lo-NeRF since the baselines are much larger than light field data.

For each ray (u, v, s, t) , the ray-based field outputs a θ_{ray} , which is then used for guiding the sampling of N points with $\{\theta_i\}_{i=0}^N$ on the ray. Each θ_i is uniformly sampled within a range dynamically adjusted along with training. On light field data, the ray-based field and the point-based field are jointly optimized. A consistent loss is adopted to update the ray-based field. The difference is that the ray-based field can be directly updated along with Hi-NeRF since the training is more stable with sparse observations on forward-facing scenes. Joint training the two fields is implemented by progressively regularizing the high-frequency radiance field, making the framework more compact and efficient. Mathematically, we set $\theta_i \sim U(\theta_{\text{ray}} - \alpha(\theta_f - \theta_n), \theta_{\text{ray}} + \alpha(\theta_f - \theta_n))$, where $\alpha \in [k, 1]$ and k is the defined range for sampling. During training, α is 1 at the beginning and then linearly converge to k . In this way, the geometry of the neural radiance field will gradually converge to local details.

Besides the reconstruction loss for supervising the radiance field, a consistent loss is adopted to update the neural light field. The consistent loss is defined as the difference between θ from the neural light field and the neural radiance field, that is, $L_{\text{consist}} = \|\theta_{\text{nelf}} - \theta_{\text{nerf}}\|_2^2$. Note that the neural light field takes low-frequency inputs, aiming for regularizing the geometry of the high-frequency radiance field. A neural light field with low-frequency inputs leads to similar effects as a neural radiance field: lacking details but generalizing well.

B. Implementation details

Our method is implemented with PyTorch [35]. Each model is trained on one 2080Ti with batch size 3072. Important hyperparameters are listed below for each experiment and other detailed training hyperparameters can be find in our code.

2D toy demonstration. The low-frequency neural field is with $L = 5$ and $s = 5$. The high-frequency neural field is with $L = 10$ and $s = 5$.

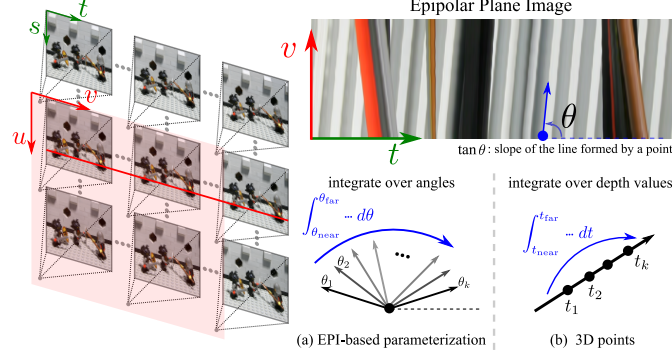


Figure 11. Illustration of Epipolar Plane Images (EPIs). Each point in the space forms a line on EPIs. To accumulate the color of a ray, integration is taken over (a) the angle θ for the EPI-based coordinates and (b) depth t for point-based coordinates.

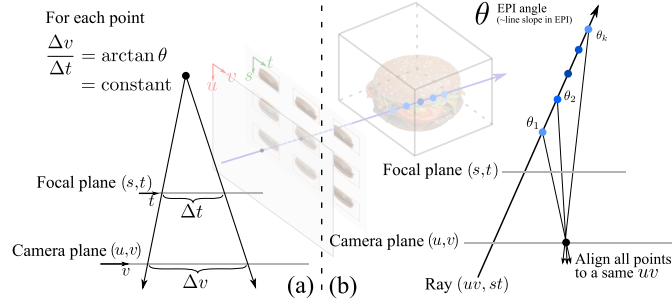


Figure 12. Points are parameterized by (s, t, θ) . (a) Once the two planes are determined, $\frac{\Delta v}{\Delta t}$ is a constant. (b) We align all rays to the same uv for representing points with 3D vectors.

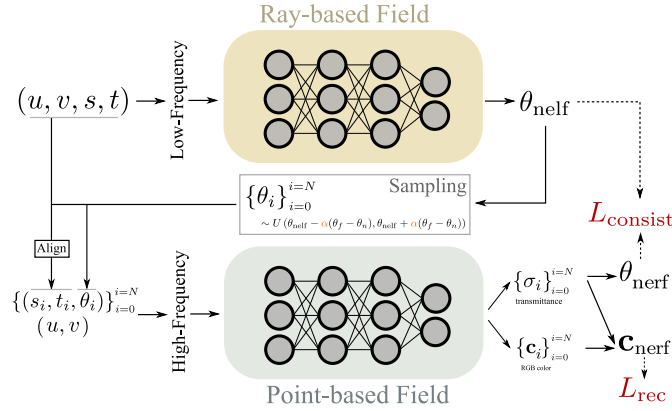


Figure 13. The overall framework of harnessing low-frequency neural fields. Inputs from 4D light field data are used for demonstration.

Experiments with 8 views. Training images for the 8 view setting are included in Tab. 6. For the low-frequency P.E., we set $L = 5$ and an extra scale parameter $s = 32$. The scale parameter indicates that all inputs are divided by s . The low-frequency radiance field is trained for 8,000 iterations and the high-frequency radiance field is trained for 120,000 iterations.

Experiments with 14 one side views. Training images are

`['r_58.png', 'r_5.png', 'r_2.png', 'r_8.png', 'r_9.png',
'r_10.png', 'r_16.png', 'r_34.png', 'r_35.png', 'r_40.png',
'r_52.png', 'r_53.png', 'r_54.png', 'r_60.png']`.

We set $s = 1$ and $L = 5$ for the low-frequency neural fields.

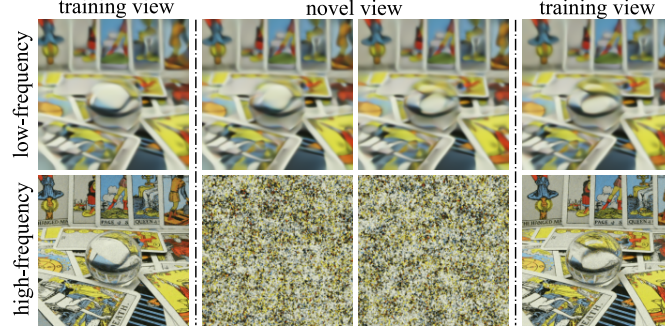


Figure 14. Rendering results of neural light fields with low- and high-frequency inputs on two adjacent training views and novel interpolation views. Using low-frequency leads to smooth color change across rays while high-frequency leads to overfitting.

Scene	Train	Val
flower	['007.png', '008.png', '011.png', '009.png']	['010.png']
leaves	['004.png', '005.png', '008.png', '006.png']	['007.png']
orchids	['008.png', '006.png', '014.png', '015.png']	['007.png']
fern	['005.png', '013.png', '015.png', '016.png']	['014.png']

Table 5. Four training views and one validation view on LLFF.

Scene	Images
lego	['r_2.png', 'r_16.png', 'r_93.png', 'r_55.png', 'r_73.png', 'r_86.png', 'r_26.png', 'r_75.png']
chair	['r_86.png', 'r_73.png', 'r_26.png', 'r_2.png', 'r_55.png', 'r_93.png', 'r_16.png', 'r_75.png']
drums	['r_86.png', 'r_93.png', 'r_75.png', 'r_26.png', 'r_55.png', 'r_73.png', 'r_16.png', 'r_2.png']
figus	['r_2.png', 'r_93.png', 'r_73.png', 'r_86.png', 'r_75.png', 'r_26.png', 'r_55.png', 'r_16.png']
mic	['r_55.png', 'r_2.png', 'r_93.png', 'r_75.png', 'r_16.png', 'r_26.png', 'r_86.png', 'r_73.png']
ship	['r_55.png', 'r_93.png', 'r_26.png', 'r_75.png', 'r_16.png', 'r_33.png', 'r_73.png', 'r_86.png']
materials	['r_75.png', 'r_26.png', 'r_93.png', 'r_55.png', 'r_86.png', 'r_16.png', 'r_2.png', 'r_73.png']
hotdog	['r_16.png', 'r_93.png', 'r_75.png', 'r_86.png', 'r_2.png', 'r_55.png', 'r_73.png', 'r_26.png']

Table 6. Training images for the 8 view setting.

Experiments on light field data. For light field data, we use the Normal distribution based P.E. proposed by [45]. The std for uv and st are 64 with length 10, so uv is the low-frequency part. The std for θ is 8 with length 5. The parameter α linearly decay to 0.5 after 20,000 iterations.

Experiments on LLFF [28]. The four training views and one validation view on the four selected scenes are as in Tab. 5.

C. Additional results

First, detailed results on each scene are included in Tab. 7 (for Tab 1 in the main text), Tab. 8 (for Tab 3 in the main text) and Tab. 9 (for Tab 3 in the main text).

Two ablation studies are conducted. Fig. 15 demonstrates the impact of different frequency settings. As we admitted in the limitation section, the best frequency for a smooth geometry needs to be determined empirically. Fortunately, as demonstrated by Fig. 15, it not hard to find an appropriate frequency for smooth effects. Also, in Fig. 16, we show the impact of the number of input views. First k views demonstrated in Tab. 7 are used for the k view setting. It can be observed that two views are not able to generate a good rough geometry, which is another limitation pointed out in the main text. Still, it is interesting to observe that 3 views are enough for a good rough geometry, though some areas are filled with undesired white points.

Per-scene performance on the dynamic scenes (D-NeRF dataset [36]) is presented in Tab. 10.

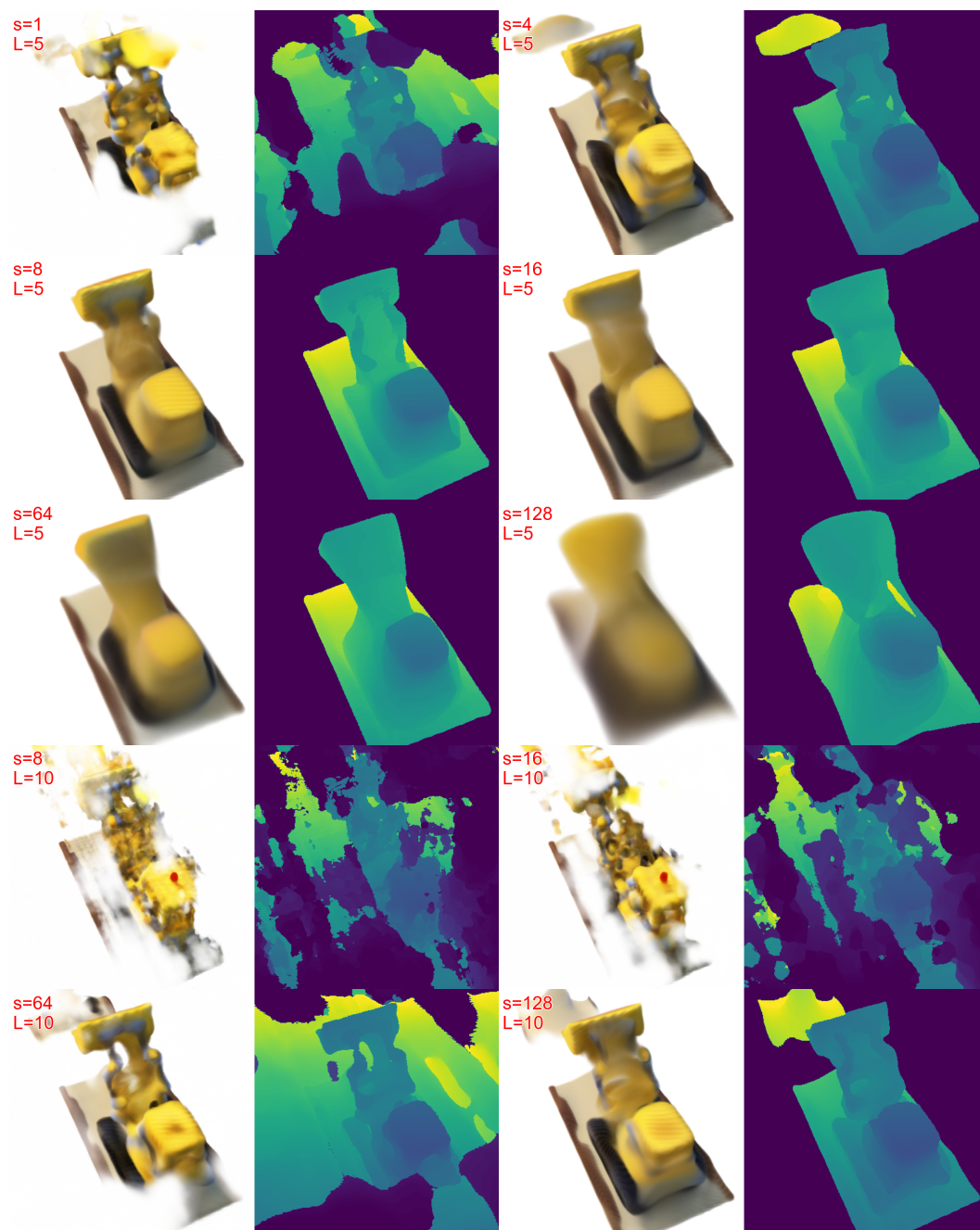


Figure 15. Results with different frequency settings (higher s and L indicates a higher frequency).

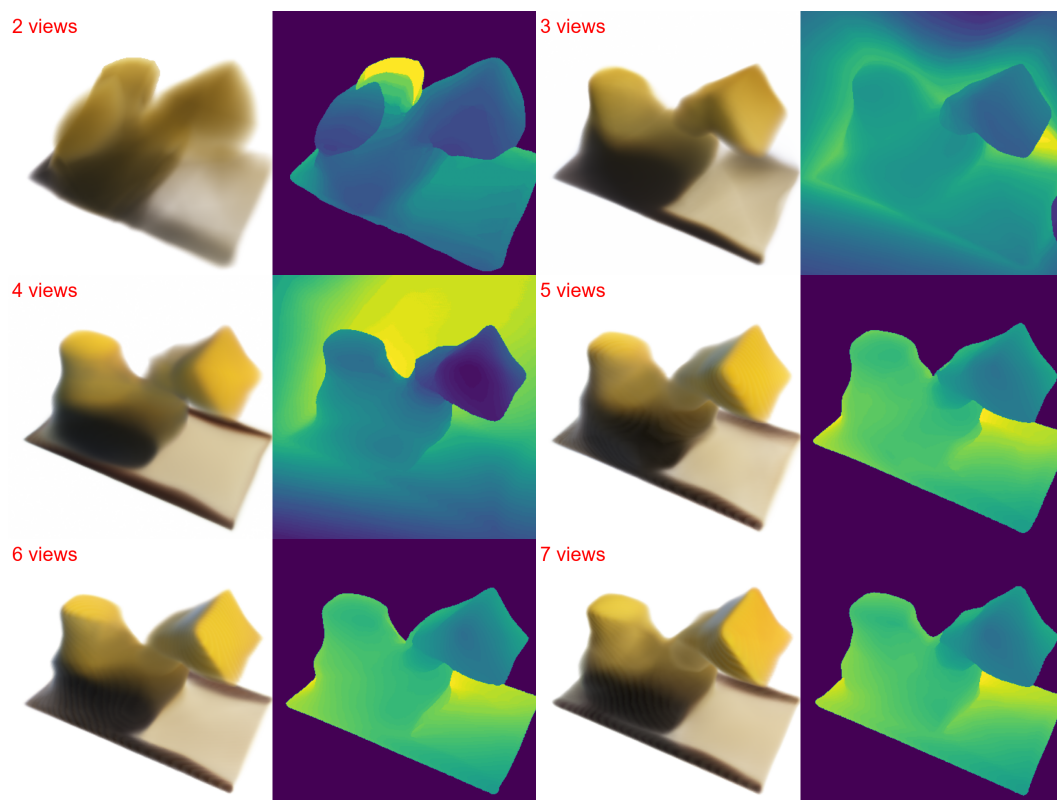


Figure 16. Results with different number of views.

PSNR \uparrow	Lego	Chair	Drums	Ficus	Mic	Ship	Materials	Hotdog
NeRF	9.726	21.049	17.472	13.728	26.287	12.929	7.837	10.446
NV	17.652	20.515	16.271	19.448	18.323	14.457	16.846	19.361
Simplified NeRF	16.735	21.870	15.021	21.091	24.206	17.092	20.659	24.060
DietNeRF	23.897	24.633	20.034	20.744	26.321	23.043	21.254	25.250
DietNeRF ft	24.311	25.595	20.029	20.940	26.794	22.536	21.621	26.626
HALO (<i>Ours</i>)	23.927	25.203	19.543	21.531	25.650	22.223	21.163	26.909
NeRF, 100 views	31.618	34.073	25.530	29.163	33.197	29.407	29.340	36.899

SSIM \uparrow	Lego	Chair	Drums	Ficus	Mic	Ship	Materials	Hotdog
NeRF	0.526	0.861	0.770	0.661	0.944	0.605	0.484	0.644
NV	0.707	0.795	0.675	0.815	0.816	0.602	0.721	0.796
Simplified NeRF	0.775	0.859	0.727	0.872	0.930	0.694	0.823	0.894
DietNeRF	0.863	0.898	0.843	0.872	0.944	0.758	0.843	0.904
DietNeRF ft	0.875	0.912	0.845	0.874	0.950	0.757	0.851	0.924
HALO (<i>Ours</i>)	0.854	0.898	0.818	0.888	0.936	0.756	0.826	0.930
NeRF, 100 views	0.965	0.978	0.929	0.966	0.979	0.875	0.958	0.981

LPIPS \downarrow	Lego	Chair	Drums	Ficus	Mic	Ship	Materials	Hotdog
NeRF	0.467	0.163	0.231	0.354	0.067	0.375	0.467	0.422
NV	0.253	0.175	0.299	0.156	0.193	0.456	0.223	0.203
Simplified NeRF	0.218	0.152	0.280	0.132	0.080	0.283	0.151	0.139
DietNeRF	0.110	0.092	0.117	0.097	0.053	0.204	0.102	0.097
DietNeRF ft	0.096	0.077	0.117	0.094	0.043	0.193	0.095	0.067
HALO (<i>Ours</i>)	0.140	0.108	0.190	0.140	0.088	0.259	0.202	0.091
NeRF, 100 views	0.033	0.025	0.064	0.035	0.023	0.125	0.037	0.025

Table 7. Detailed results on each scene for the 8 views setting.

Scene	NeRF			HALO		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
amethyst	23.146	0.708	0.342	27.088	0.797	0.262
bracelet	22.215	0.794	0.299	28.929	0.936	0.188
cards-big	17.485	0.535	0.507	20.613	0.778	0.260
cards-small	23.199	0.868	0.163	28.382	0.938	0.102
chess	27.707	0.858	0.337	33.682	0.946	0.227
eucalyptus-flowers	35.491	0.930	0.277	38.163	0.950	0.233
jellybeans	35.297	0.974	0.169	39.584	0.982	0.153
lego-bulldozer	24.700	0.774	0.482	28.688	0.894	0.330
lego-gantry	17.809	0.625	0.433	19.741	0.672	0.523
lego-knights	27.210	0.899	0.170	32.494	0.950	0.112
lego-truck	35.099	0.945	0.256	38.053	0.960	0.224
stanfordbunny	39.046	0.945	0.215	42.146	0.967	0.198
treasure	14.671	0.218	0.511	29.127	0.895	0.242
mean	26.390	0.775	0.320	31.283	0.897	0.234

Table 8. Detailed results on each scene for the light field data StanfordLF.

	PSNR \uparrow					SSIM \uparrow					LPIPS \downarrow				
	NeRF	IBRNet	MVSNeRF	HALO w/o j.t	HALO	NeRF	IBRNet	MVSNeRF	HALO w/o j.t	HALO	NeRF	IBRNet	MVSNeRF	HALO w/o j.t	HALO
fern	23.995	23.225	23.723	22.64	23.10	0.763	0.716	0.727	0.774	0.795	0.274	0.302	0.299	0.266	0.253
flower	22.247	23.052	23.372	26.55	27.23	0.739	0.698	0.719	0.909	0.912	0.282	0.251	0.233	0.146	0.143
leaves	18.542	18.365	18.694	22.07	21.54	0.652	0.561	0.603	0.843	0.826	0.296	0.371	0.331	0.180	0.222
orchids	14.844	16.133	16.522	19.01	20.51	0.383	0.343	0.458	0.705	0.732	0.508	0.559	0.425	0.286	0.258
mean	19.907	20.194	20.578	22.567	23.095	0.634	0.580	0.627	0.808	0.816	0.340	0.371	0.322	0.219	0.219

Table 9. Detailed results on each scene for the 4 views setting on LLFF.

Scene	PSNR \uparrow		SSIM \uparrow		LPIPS \downarrow	
	TNV[9]	+ Ours	TNV[9]	+ Ours	TNV[9]	+ Ours
HELLWARRIOR	12.288	17.558	0.733	0.832	0.349	0.247
BOUNCINGBALLS	25.431	24.168	0.942	0.911	0.135	0.116
JUMPINGJACKS	22.591	31.952	0.915	0.972	0.114	0.087
HOOK	16.676	27.065	0.692	0.953	0.370	0.064
LEGO	19.824	20.083	0.845	0.787	0.203	0.244
MUTANT	26.154	28.953	0.935	0.952	0.077	0.054
STANDUP	24.404	28.980	0.934	0.963	0.089	0.042
TREX	23.770	26.736	0.912	0.939	0.120	0.076
Average	21.392	25.687	0.864	0.914	0.182	0.116

Table 10. Few-shot rendering on dynamic scenes. Every 8 training images (*i.e.*, sparse spatiotemporal inputs) from the D-NeRF [36] dataset are used.