

Please cite this paper as:

M. Bonotto, L. Sarrocco, D. Evangelista, M. Imperoli, and A. Pretto "CombiNeRF: A Combination of Regularization Techniques for Few-Shot Neural Radiance Field View Synthesis," in International Conference on 3D Vision (3DV), 2024.

CombiNeRF: A Combination of Regularization Techniques for Few-Shot Neural Radiance Field View Synthesis

Matteo Bonotto^{1,2*} Luigi Sarrocco^{1,2*} Daniele Evangelista^{1,2} Marco Imperoli² Alberto Pretto¹¹Department of Information Engineering, University of Padova, Italy²FlexSight, Via Prima Strada 35, Padova, Italy

Abstract

Neural Radiance Fields (NeRFs) have shown impressive results for novel view synthesis when a sufficiently large amount of views are available. When dealing with few-shot settings, i.e. with a small set of input views, the training could overfit those views, leading to artifacts and geometric and chromatic inconsistencies in the resulting rendering. Regularization is a valid solution that helps NeRF generalization. On the other hand, each of the most recent NeRF regularization techniques aim to mitigate a specific rendering problem. Starting from this observation, in this paper we propose CombiNeRF, a framework that synergically combines several regularization techniques, some of them novel, in order to unify the benefits of each. In particular, we regularize single and neighboring rays distributions and we add a smoothness term to regularize near geometries. After these geometric approaches, we propose to exploit Lipschitz regularization to both NeRF density and color networks and to use encoding masks for input features regularization. We show that CombiNeRF outperforms the state-of-the-art methods with few-shot settings in several publicly available datasets. We also present an ablation study on the LLFF and NeRF-Synthetic datasets that support the choices made. We release with this paper the open-source implementation of our framework.

1. Introduction

Neural Radiance Field [16] has emerged as a powerful approach for scene reconstruction and photorealistic rendering from a sparse set of 2D images. By leveraging multi-layer perceptron (MLP) to model the volumetric scene representation, NeRF can generate high-quality novel views of scenes. As the MLP networks are queried multiple times for

*Authors contributed equally to this work. This research is part of a FlexSight-supported project that has received funding from the European Union's Horizon 2020 research and innovation programme – IN-NOSUP under grant agreement No 101005711. Corresponding author: matteo.bonotto.2@phd.unipd.it

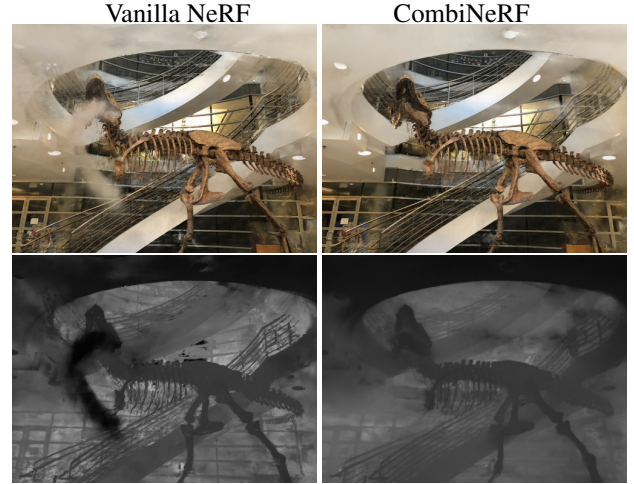


Figure 1. The figure shows how the proposed CombiNeRF achieves better results in terms of rendering and reconstruction quality in few-shot settings compared with the Vanilla NeRF [18, 34].

each pixel to be rendered, a lightweight architecture allows to dramatically speed up the whole rendering pipeline. The Multi-resolution Hash Encoding proposed in [18] allows, for example, to improve scene reconstruction efficiency using a shallow MLP without any loss of visual accuracy. However, to achieve high accuracy in its predictions and avoid artifacts (see Fig. 1), NeRF often relies on a large number of images, which can be impractical in real-world scenarios.

Regularization is a crucial tool for improving the visual fidelity of rendered images, ensuring more coherent predictions. Several works [3, 10, 23, 42] have proposed to constrain the MLP network during training by adding regularization terms in the final loss. Instead, Lipschitz regularization [13, 29] is applied directly to network weights to enforce smoothness by controlling the rate of change of the network's output concerning its inputs.

In this work, we empirically select a combination of loss components to constrain the MLP network during the training phase. Firstly, (a) we modify the information-theoretic

approach of InfoNeRF [10] that includes regularization of neighboring rays distributions, (b) we adopt a smoothness term to regularize near geometries, and (c) we regularize single ray distributions in order to have high density values in regions corresponding to object surface. Additionally, as in [11], we add an encoding mask that encourages the learning of finer details only in the latest stages of the training phase. Finally, we impose Lipschitz regularization on both color and density networks, observing substantial improvements in the overall reconstruction and rendering quality and showing how NeRF benefits from Lipschitz regularization. CombiNeRF avoids the need for pre-training required by similar approaches while showing promising improvements over the state of the art.

In summary, we present the following contributions: *i)* we propose a modification of the KL-Divergence loss proposed in InfoNeRF; *ii)* we select a combination of regularization losses, available in a unified framework, and we validate its effectiveness through ablation studies; *iii)* to our knowledge we are the first to impose Lipschitz regularization on all the network layers on NeRF, recording a performance increase in all the tested scenarios; *iv)* an extensive performance evaluation on public datasets, showing that CombiNeRF achieves state-of-the-art (SOTA) performance in few-shot setting; *v)* we open-sourced the CombiNeRF implementation¹.

2. Related Work

Novel view synthesis aims to generate images from different viewpoints by utilizing a collection of pre-existing views [4, 14, 20, 30, 31]. Among the novel view synthesis solutions, NeRF emerged as a result of several factors as the compactness offered by the underlying structure, its domain-agnostic nature and the impressive visual quality offered. Subsequent works managed to improve over the fidelity of rendered images [1, 17, 33, 35, 41], reduce required time [19, 24, 28, 38] and extend NeRF capabilities and application domain [7, 21, 25, 26, 43, 51].

NeRF Regularization. NeRF models struggle to render novel high-quality images when supervised with only a few input views during training. To reduce the impact of the artifacts introduced by overfitting on training samples, previous works [2, 6, 8, 23, 32, 46] focused on adding loss terms to provide additional constraints to the model, integrating prior knowledge and proposing a custom pipeline. In InfoNeRF [10] the entropy of the rays and the KL-divergence between the distribution of neighboring rays are both minimized in order to restrict the range of the prediction of high-density values on the object surface. Similarly, in Mip-NeRF 360 [3] distortion

loss was introduced to consolidate weights into as small a region as possible when rays intersect a scene object. In RegNeRF and DiffusioNeRF [23, 42] rendered patches that are less likely to appear, according to a determined image distribution, are penalized. In PixelNeRF [46] prior knowledge acquired from different scenes and global context information are both leveraged by concatenating image features extracted by a Convolutional Neural Network Encoder with the input 3D positions. In Aug-NeRF [5] noise is injected into input, feature and output during training to improve the generalization capabilities.

Neural Surface Reconstruction. Implicit functions like signed distance functions (SDFs) [36, 45, 47, 48] and occupancy maps [12, 22] are best suited for representing objects on the scene with a defined surface geometry. Both NeuS and HF-NeuS [37, 39] reparametrize the rendering equation used in NeRF to exploit SDFs properties. In order to increase the training and inference speed and facilitate the learning of high-frequency details, Instant-NSR and Neuralangelo [11, 50] exploit the hash grid encoding proposed in Instant-NGP [18], while PermutoSDF [27] employs a permutohedral lattice to decrease the memory accesses. To recover smoother surfaces in reflective or untextured areas, Neuralangelo and PermutoSDF add a curvature term loss. Additionally, Neuralangelo discards values from the finer levels of the hash grid encoding to encourage the learning of coarse details on the first training iterations. To prevent the color network overfitting due to the injection of the curvature loss term, PermutoSDF employs the Lipschitz constant regularization method proposed by Liu *et al.* [13].

In our work, we exploit both Lipschitz regularization, leveraged in neural surface reconstruction to avoid geometry over-smoothness, and an encoding mask, finding out they provide additional benefits in the few-shot setting.

3. Method

A graphical overview of the CombiNeRF method proposed in this work is given in Fig. 2. The proposed approach combines multiple regularization techniques in a unified and flexible framework that helps NeRF generalization in few-shot scenarios. In this section, after providing a preliminary theoretic background, all the combined losses will be described together with the Encoding Mask and Lipschitz regularization techniques. For each technique, our improvements and modifications will be described in detail to better highlight the main contributions of this work.

3.1. Preliminaries

Neural Radiance Fields. Neural Radiance Field [16] exploits an MLP to encode the scene. For each input 3D point and view direction, it outputs the density σ of the point and

¹<https://github.com/SarroccoLuigi/CombiNeRF>

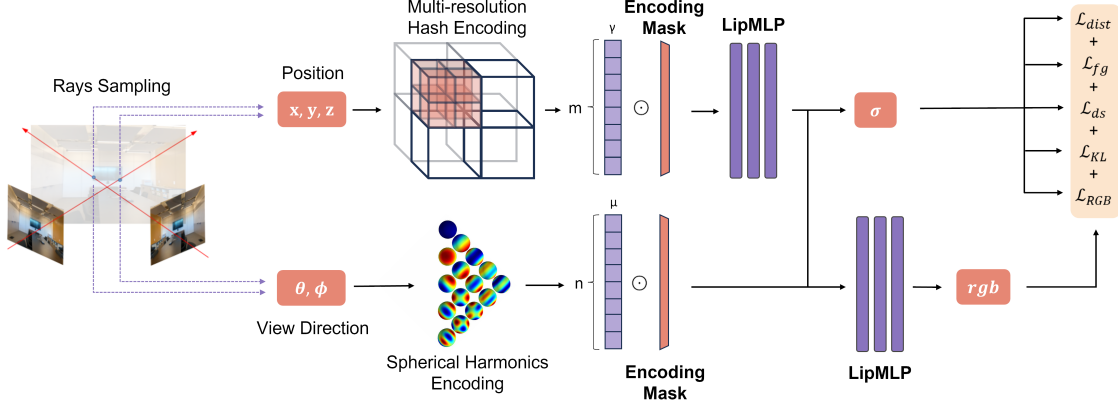


Figure 2. Overview of the CombiNeRF framework. We sample 3D points over a batch of rays passing through the scene. Position and view direction are respectively encoded through Multi-resolution Hash Encoding and Spherical Harmonics and fed to the Lipschitz network (LipMLP) after being masked. Networks’ outputs are used by volumetric rendering for estimating the expected color C and depth d of each ray, while different loss terms are computed to regularize the training process. CombiNeRF combines *i*) all these regularization losses, *ii*) the Lipschitz network instead of the original MLP, *iii*) the Encoding Mask approach used for masking the networks’ input.

its color value c . For each pixel, points are sampled on its corresponding camera ray $r(\mathbf{t}) = o + \mathbf{t}d$ and fed to the MLP network. The rendered color is approximated by:

$$C(r) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) c_i, \quad (1)$$

where $\delta_i = \mathbf{t}_{i+1} - \mathbf{t}_i$ is the distance between the i^{th} point and its adjacent sample, N is the number of samples and the transmittance T_i is computed as:

$$T_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right). \quad (2)$$

The rendered pixel color $C(r)$ can be expressed in terms of a weighted sum of the color values $C(r) = \sum_{i=1}^N w_i c_i$ where $w_i = T_i \alpha_i$ with $\alpha_i = (1 - \exp(-\sigma_i \delta_i))$. The final loss to train the network is the total squared error between the rendered and true pixel colors:

$$\mathcal{L}_{RGB} = \sum_{r \in \mathcal{R}} \|C(r) - C^*(r)\|_2^2, \quad (3)$$

where \mathcal{R} is the set of rays in each batch, and $C^*(r)$ is the ground truth RGB color for ray r .

Multi-resolution Hash Encoding. In NeRF, a sine-based positional encoding is used to map the 3D point position to a high-dimensional feature vector in order to capture high-frequency details. On the other hand, the Multi-resolution Hash Encoding proposed by Müller *et al.* [18] adopts multi-resolution grids, achieving a similar result while significantly accelerating training and inference times. This is achieved by first identifying the grid cell

containing the point and then assigning the result of the spatial interpolation of the features defined on the corners. As a result, a lightweight MLP can be employed to process the encoded points. Given an input point $x \in \mathbb{R}^3$, its feature vector $\gamma(x) \in \mathbb{R}^{FL}$ can be expressed as:

$$\gamma(x) = [\gamma_{1,1}(x), \dots, \gamma_{F,1}(x), \dots, \gamma_{F,L}(x)], \quad (4)$$

where F is the dimension of the feature space and L is the number of levels. To minimize memory usage, each grid corner index is mapped into a hash table that holds its corresponding feature vector. Collisions are limited to finer levels where $(N_l + 1)^3 > E$, N_l and E being the resolution at level l and the maximum number of entries in the hash maps respectively, and they are handled through gradient averaging. This ensures that the results maintain comparable quality to using a multi-resolution grid without hash encoding.

3.2. Regularization Losses

KL-Divergence Loss. When dealing with a few training views, the standard models [16, 18] struggle to generalize to unseen views. As 3D scenes exhibit piece-wise smooth surfaces, we enforce similar distributions of weight values among neighboring rays. Differently from [10], we compute the ray density $p(\mathbf{r})$ on weights w_i , instead of alpha values α_i :

$$p(\mathbf{r}_i) = \frac{w_i}{\sum_j w_j} = \frac{T_i \alpha_i}{\sum_j T_j \alpha_j}. \quad (5)$$

Given an observed ray \mathbf{r} , we sample another neighboring ray $\hat{\mathbf{r}}$ minimizing the KL-Divergence (D_{KL}) between their weight distribution. The corresponding loss is defined as

follows:

$$\mathcal{L}_{KL} = D_{KL}\left(P(\mathbf{r})||P(\hat{\mathbf{r}})\right) = \sum_{i=1}^N p(\mathbf{r}_i) \log \frac{p(\mathbf{r}_i)}{p(\hat{\mathbf{r}}_i)}. \quad (6)$$

We select the neighboring ray by choosing one of the possible adjacent pixel rays. After sampling uniformly at random one of the four adjacent rays for each training ray, we compute \mathcal{L}_{KL} .

Distortion and Full Geometry Loss. A common problem with NeRF is the presence of "floaters", disconnected regions of dense space, usually located close to the camera planes, which cause the presence of blurry cloud artifacts on images rendered from novel views. Previous works [3, 42] showed that those artifacts can be removed by adding a loss term \mathcal{L}_{dist} that takes into account the distances t_i and weights w_i of the N points sampled on the rays, and the depth

$$d(r) = \frac{\sum_{i=1}^N w_i t_i}{\sum_{i=1}^N w_i}, \quad (7)$$

of each ray. Thus, the \mathcal{L}_{dist} loss can be written as:

$$\mathcal{L}_{dist} = \frac{1}{d(r)} \left(\sum_{i,j} w_i w_j \left| \frac{t_i + t_{i+1}}{2} - \frac{t_j + t_{j+1}}{2} \right| + \frac{1}{3} \sum_{i=1}^N w_i^2 (t_{i+1} - t_i) \right), \quad (8)$$

where the depth factor penalizes dense regions near the camera. The first term minimizes the weighted distances between all pairs of interval midpoints and the second term minimizes the weighted size of each individual interval. In this way, weights on the ray are encouraged to be as compact as possible by pulling distant intervals towards each other, consolidating weights into a single interval or a small number of nearby intervals, and minimizing the width of each interval.

In Eq. (8), all possible intervals for each ray are considered. As the total combination of all (i, j) would drastically increase the amount of memory required, to reduce the total load we consider only a subset over the total number of sampled rays on which we compute the loss.

Furthermore, a normalization loss term \mathcal{L}_{fg} is considered, following [42], that encourages the weights to sum to unit, as the ray is expected to be fully absorbed by the scene geometry in real scenarios:

$$\mathcal{L}_{fg} = \left(1 - \sum_{i=1}^N w_i \right)^2, \quad (9)$$

In addition \mathcal{L}_{fg} enforces \mathcal{L}_{KL} by driving the model to treat the weights w_i as probabilities, $p(\mathbf{r}_i) \approx w_i$ (See

Eq. (5)).

Depth Smoothness Loss. Similarly to KL-Divergence loss, the depth smoothness constraint encourages smooth surfaces. However, while the former term requires similarity between the compared distributions, the latter term only computes the difference between the expected depth values. As in [23], given the estimated depth $d(r)$ of a pixel ray and a patch of size S_{patch} , the depth smoothness constraint is defined as follows:

$$\mathcal{L}_{ds} = \sum_{p \in P} \sum_{i,j=1}^{S_{patch}-1} \left((d(r_{ij}) - d(r_{i+1j}))^2 + (d(r_{ij}) - d(r_{ij+1}))^2 \right), \quad (10)$$

where P is the set of all the rendered patches and r_{ij} is the ray passing through the pixel (i, j) of patch p .

3.3. Encoding Mask

Input encoding allows preserving of high-frequency details [16]. However, when only few images are available, the network is more sensitive to noise. In this scenario, high-frequency components exacerbate the problem, preventing the network from exploring more in depth low-frequency information and consequently learning a coherent geometry. As in [11, 44], we use a mask to filter out high-frequency components of the input in early iterations in order to let the network focus on robust low-frequency information. Given the length $l = L \cdot F$ of the resulting Multi-resolution Hash Encoding given by Eq. (4), the mask m is defined as:

$$m = [1_1, \dots, 1_{\lfloor l \cdot x \rfloor}, 0, \dots, 0_l], \quad (11)$$

where x is the ratio of features to keep active. The resulting feature vector given in input to the MLP network will be the element-wise product $\gamma \odot m$ between the Multi-resolution Hash Encoding of the input and the mask. The ratio x is set to keep only the coarsest features during the initial phase of the training, progressively including the remaining features. In CombiNeRF, we apply Instant-NGP Multi-resolution Hash Encoding on input positions and Spherical Harmonics Encoding on the view directions (see Fig. 2). The progressive mask can be used on both inputs in order to filter high-frequency terms on early training iterations.

3.4. Lipschitz Network

We define $f_{\Theta}(x, t)$ as the function representing the implicit shape, expressed by means of a neural network, where $\Theta = \{W_i, b_i\}$ is the set of parameters containing weights W_i and biases b_i of each layer l_i of the network.

A function is called Lipschitz continuous if there exists a constant $k \geq 0$ such that:

$$\|f_{\Theta}(x_0) - f_{\Theta}(x_1)\| \leq k \|x_0 - x_1\|, \quad (12)$$

for all possible inputs x_0, x_1 , where k is called the Lipschitz constant which bounds how fast the function f_Θ can change.

We employ the regularization method proposed in [13], which was already used by PermutoSDF [27] on the color network to balance out the effect introduced by the curvature regularization applied to the SDF network. However, differently from PermutoSDF we apply the regularization to both NeRF density and color networks, recording an increase in performance with respect to the original MLP network. Given the trainable Lipschitz bound k_i associated to layer l_i , the normalization scheme is described as:

$$\begin{aligned} y &= \sigma(\hat{W}_i x + b_i) \\ \hat{W}_i &= f_n(W_i, \ln(1 + e^{k_i})), \end{aligned} \quad (13)$$

where f_n is a normalization function and the term $\ln(1 + e^{k_i})$ allows to avoid negative values for k_i . The normalization scales each row of W_i to have an absolute value row-sum less or equal to $\ln(1 + e^{k_i})$ in order to satisfy Eq. (12).

3.5. Overall Method

CombiNeRF combines the previously described regularization techniques regarding losses and network structure, hence the name *CombiNeRF*. Thus, we can write the final loss as:

$$\begin{aligned} \mathcal{L}_{CombiNeRF} &= \mathcal{L}_{RGB} + \lambda_{dist} \cdot \mathcal{L}_{dist} + \\ &\lambda_{fg} \cdot \mathcal{L}_{fg} + \lambda_{ds} \cdot \mathcal{L}_{ds} + \lambda_{KL} \cdot \mathcal{L}_{KL}, \end{aligned} \quad (14)$$

where λ are the hyperparameters controlling the contribution of each loss. In addition, CombiNeRF includes Lipschitz regularization and the Encoding Mask technique. The proposed CombiNeRF offers a unified implementation of all the regularization techniques described above, outperforming current SOTA methods on few-shot scenarios as demonstrated in the following experimental section.

4. Experiments

In this section we present the dataset and metrics we use for the experiments, we provide relevant details about CombiNeRF and we show quantitative and qualitative comparisons made against the SOTA methods. Finally, an ablation study is made to evaluate the contribution of each component of which CombiNeRF is built.

4.1. Setup

Dataset. Our experiments include the LLFF dataset [15] and the NeRF-Synthetic dataset [16] under few-shot settings. LLFF is composed by 8 complex scenes, representing real-world scenarios, with 20-62 images for each scene captured with a consumer camera. Instead, NeRF-Synthetic is composed of 8 synthetic scenes with view-dependent light

transport effects, in which each scene is composed of 100 training images and 200 test images.

We evaluate LLFF on 3/6/9 input views, following the protocol proposed in RegNeRF [23], while in NeRF-Synthetic we train 8 views and test 25 views following DietNeRF [9].

Metrics. To evaluate the performance of CombiNeRF against the SOTA methods, we use 3 different metrics: peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM) [40], and learned perceptual image patch similarity (LPIPS) [49]. We also take into consideration the geometric mean of $\text{MSE} = 10^{-\text{PSNR}/10}$, $\sqrt{1 - \text{SSIM}}$ and LPIPS following [23], in order to have an easier and unified comparison.

Implementation Details. For our experiments, we used the torch-ngp [34] implementation of Instant-NGP [18] as base code, since from our experience it generally outperforms the original NeRF implementation [16] from both a runtime and accuracy points of view. We refer to torch-ngp as "Vanilla NeRF" in the experiments. We developed CombiNeRF on top of torch-ngp, thus obtaining a unified and unique implementation that embeds all the contributions shown in the previous sections. As already mentioned, the final implementation of CombiNeRF is made publicly available with this paper.

For the LLFF dataset, we set $\lambda_{dist} = 0$ for the first 1000 iterations and then $\lambda_{dist} = 2 \cdot 10^{-5}$ until the end and $\lambda_{fg} = 10^{-4}$, $\lambda_{KL} = 10^{-5}$, $\lambda_{ds} = 0.1$ and $S_{patch} = 4$, while we use the encoding mask only on the density network in which x saturates after 90% of the total iterations. In the NeRF-Synthetic dataset we set $\lambda_{dist} = 0$ for the first 1000 iterations and then $\lambda_{dist} = 2 \cdot 10^{-3}$ until the end and $\lambda_{fg} = 10^{-3}$, $\lambda_{KL} = 10^{-5}$, $\lambda_{ds} = 0.01$ and $S_{patch} = 4$, x saturates after 20% of the total iterations. We also set the number of sampled rays for each iteration to 4096 and 7008 and the number of levels for the Instant-NGP Multi-resolution Hash Encoding to 16 and 32 for LLFF and NeRF-Synthetic, respectively.

4.2. Comparison

We compare the performance achieved by CombiNeRF against the SOTA methods in the few-shot scenario. In LLFF we consider RegNeRF [23], FreeNeRF [44] and DiffusioNeRF [42], while in NeRF-Synthetic we take into account DietNeRF [9] and FreeNeRF [44]. All quantitative evaluation results for the other methods, along with images showing the qualitative results, are taken from the respective papers. If some results are not present, it means that the related paper has not reported quantitative or qualitative results on the related dataset.

	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow			Average \downarrow		
	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
Vanilla NeRF	17.71	22.03	24.21	0.544	0.738	0.811	0.303	0.149	0.101	0.155	0.081	0.057
RegNeRF	19.08	23.10	24.86	0.587	0.760	0.820	0.336	0.206	0.161	0.146	0.086	0.067
FreeNeRF	19.63	23.73	25.13	0.612	0.779	0.827	0.308	0.195	0.160	0.134	0.075	0.064
DiffusioNeRF	19.88	24.28	25.10	0.590	0.765	0.802	0.192	0.101	0.084	0.118	0.071	0.060
CombiNeRF	20.37	23.99	25.15	0.686	0.805	0.841	0.191	0.106	0.084	0.101	0.060	0.049

Table 1. Comparison of CombiNeRF with SOTA methods on the LLFF dataset with 3/6/9 input view few-shot settings. DiffusioNeRF row relates to the "Geometric Baseline" of the corresponding paper, as it has the highest metric scores.

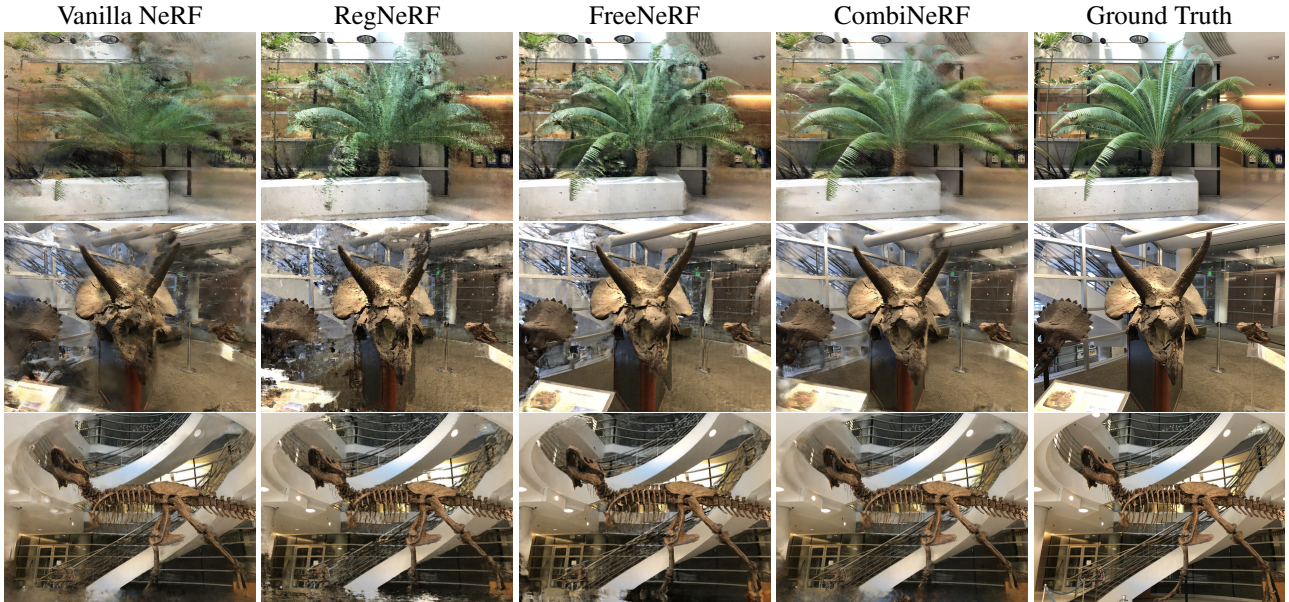


Figure 3. Comparison of our CombiNeRF against RegNeRF, FreeNeRF and Vanilla NeRF on Fern, Horns and Trex scenarios with 3-view setting.

LLFF Dataset. This dataset is a collection of complex real-world scenes where, from our experiments, we observed that Vanilla NeRF heavily overfits the training images when only a few of them are used as input. Tab. 1 shows quantitative results under 3/6/9 training images. For 3 and 9 views, CombiNeRF outperforms the other methods in all the metrics. With 6 views, DiffusioNeRF achieves a better LPIPS score, however, CombiNeRF scores higher on average. Interestingly, with more views (6 or 9), Vanilla NeRF approaches and sometimes outperforms all other methods except CombiNeRF, while its performance heavily degrades with a reduced number of views. CombiNeRF, on the other hand, shows cutting-edge and consistent performance regardless of the number of views.

Qualitative results on LLFF are shown in Fig. 3. CombiNeRF is able to reconstruct better high-frequency details maintaining the geometry of the scene. Other methods generate lots of artifacts, like in the "Fern" scene, and they struggle to reconstruct the background. In the "Horns" scene, the environment is noisy in RegNeRF and FreeNeRF

while it becomes sharper in CombiNeRF. We also notice that some fine details like the stair handrail in the "Trex" scene are preserved. In Fig. 4 we additionally compare CombiNeRF against FreeNeRF. In "Room" with the 9-view setting, the near part of the table gets deformed while in CombiNeRF some details on the floor are still visible. In "Orchids" with the 6-view setting, CombiNeRF is able to render more high-frequency details, for example the stripes in the orchid petals. In "Leaves" with the 3-view setting, the central leaf rendered in FreeNeRF contains several artifacts and, in general, even in this case, CombiNeRF is able to reconstruct better high-frequency details.

NeRF-Synthetic Dataset. This dataset contains synthetic renderings of objects exhibiting high-frequency geometric details and reflective effects, thus making this dataset particularly challenging. In Tab. 2 we show the quantitative results of CombiNeRF and the compared methods. We can see that CombiNeRF outperforms the other SOTA methods on the PSNR and LPIPS metrics

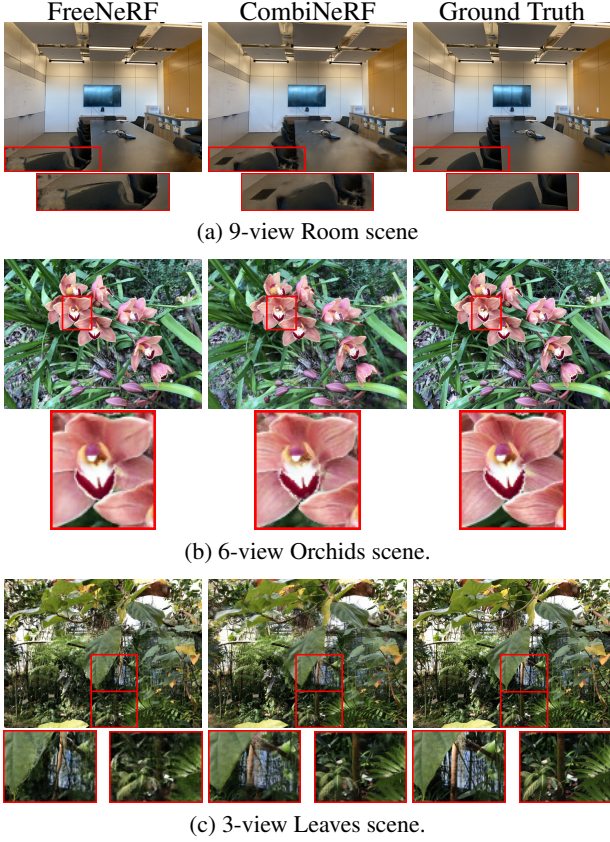


Figure 4. In-depth comparison of CombiNeRF against FreeNeRF on some LLFF scenes with 3/6/9 input views.

NeRF-Synthetic 8-views	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Vanilla NeRF	22.335	0.845	0.144
DietNeRF	23.147	0.866	0.109
DietNeRF, \mathcal{L}_{MSE} ft	23.591	0.874	0.097
FreeNeRF	24.259	0.883	0.098
CombiNeRF	24.394	0.883	0.088

Table 2. NeRF SOTA comparison on NeRF-Synthetic dataset with 8 input views. " \mathcal{L}_{MSE} ft" is the fine-tuned version of DietNeRF.

while achieving the same result on the SSIM. CombiNeRF is the overall best-performing method also on the NeRF-Synthetic dataset, performing better on scenes like "Lego", "Mic", and "Ficus" which exhibit complex geometries, while "Materials" represents the most challenging scenario due to the presence of strong view-dependent reflection.

In Fig. 5 we show qualitative results of CombiNeRF in different scenes. The reconstructed depth is more consistent, presenting fewer artifacts. Rendered images are far less noisy with respect to the ones rendered by the Vanilla NeRF implementation. Colors are better preserved, as can be seen in the "Hotdog" and "Ficus" scenes where the color of the leaves is more realistic and closer to the ground truth.

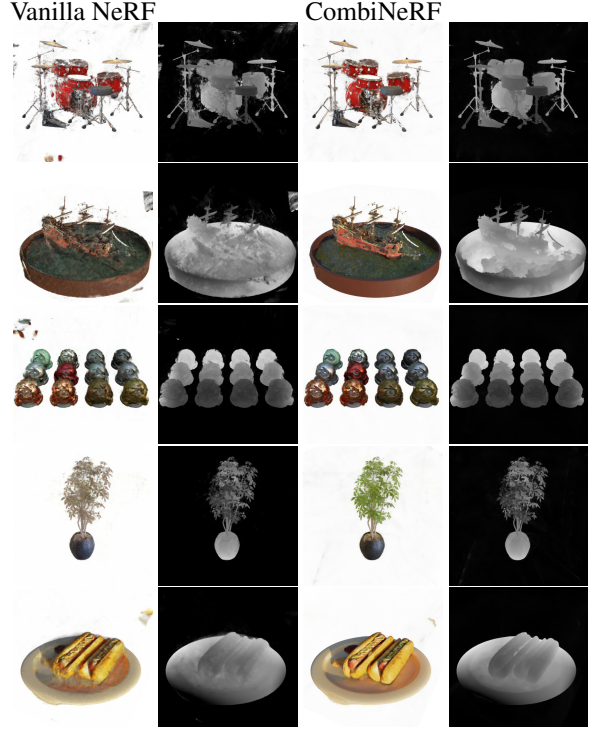


Figure 5. Comparison of CombiNeRF against the Vanilla NeRF method in Drums, Ship, Materials, Ficus and Hotdog scenarios.

4.3. Ablation Study

In this section, we assess the contribution of each of the regularization techniques used in CombiNeRF, thus highlighting their importance. Instead of showing all the possible combinations, we empirically select only the more relevant regularization techniques. We conducted an ablation study of our CombiNeRF in the 3-view and 8-view few-shot settings for the LLFF and NeRF-Synthetic datasets respectively. For this ablation study, we used the same parameters of the experiments previously described.

LLFF Dataset. In Tab. 3 we show the ablation on the LLFF dataset with the 3-view setting. Starting from the Vanilla NeRF solution, we notice that each contribution in CombiNeRF brings an increase in performance, with Lipschitz regularization and $\mathcal{L}_{dist} + \mathcal{L}_{fg}$ losses playing a crucial role in improving the overall performance. Applying Lipschitz regularization on both density and color networks also allows to further improve the quality of the obtained results. From the same table, we can observe that when using the term losses $\mathcal{L}_{dist} + \mathcal{L}_{fg}$ greatly benefit from sampling patches of rays instead of single rays.

In Fig. 6 we show qualitative results of CombiNeRF using Lipschitz regularization in both sigma and color network against the only color network. The latter tentative approach is named CombiNeRF[#] in the figures and tables.

LLFF 3-views	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
Vanilla NeRF	17.71	0.544	0.303	0.155
L	18.74	0.608	0.249	0.130
$\mathcal{L}_{dist} + \mathcal{L}_{fg}$	17.93	0.554	0.288	0.151
$\mathcal{L}_{dist} + \mathcal{L}_{fg} (S_{patch}=4)$	19.10	0.608	0.249	0.128
$\mathcal{L}_{dist} + \mathcal{L}_{fg} + \mathcal{L}_{ds}$	19.50	0.626	0.237	0.122
$\mathcal{L}_{dist} + \mathcal{L}_{fg} + \mathcal{L}_{ds} + L$	20.09	0.674	0.198	0.105
$\mathcal{L}_{dist} + \mathcal{L}_{fg} + \mathcal{L}_{ds} + L + EM$	20.27	0.683	0.194	0.103
CombiNeRF [#]	19.52	0.637	0.236	0.119
CombiNeRF	20.37	0.686	0.191	0.101

Table 3. Ablation on LLFF dataset with 3-view setting. Lipschitz and Encoding Mask are named L and EM respectively. When using \mathcal{L}_{ds} , we set the patch size $S_{patch} = 4$ by default. We call CombiNeRF[#] our method using Lipschitz only in the color network.

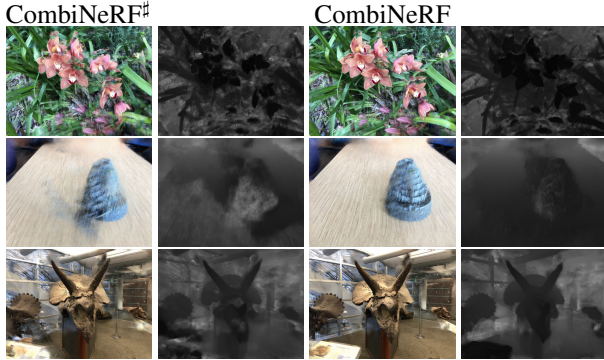


Figure 6. Qualitative result on the ablation study of LLFF with 3-view setting. We compare our CombiNeRF with CombiNeRF[#] (our method using Lipschitz only in the color network).

NeRF-Synthetic 8-views	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
Vanilla NeRF	22.34	0.845	0.144	0.073
$\mathcal{L}_{dist} + \mathcal{L}_{fg} + \mathcal{L}_{ds} + L$	23.61	0.865	0.115	0.059
$\mathcal{L}_{dist} + \mathcal{L}_{fg} + \mathcal{L}_{ds} + \mathcal{L}_{KL}^\dagger + L$	22.93	0.863	0.126	0.064
$\mathcal{L}_{dist} + \mathcal{L}_{fg} + \mathcal{L}_{ds} + L + EM$	24.04	0.871	0.105	0.056
CombiNeRF [†]	24.36	0.883	0.088	0.050
CombiNeRF	24.39	0.883	0.088	0.051

Table 4. Ablation on NeRF-Synthetic dataset with 8-view setting. Lipschitz is defined as L and Encoding Mask is defined as EM. We call \mathcal{L}_{KL}^\dagger the KL-Divergence loss presented by InfoNeRF and we call CombiNeRF[†] our method using their \mathcal{L}_{KL}^\dagger .

CombiNeRF is able to model in a more coherent way the geometry of the scenarios, removing also the blurry effects which characterize the renderings of the implementation using only Lipschitz regularization in the color network.

NeRF-Synthetic Dataset. As already done for the LLFF dataset, we also provide an ablation study on the NeRF-Synthetic dataset. In Tab. 4 we show the result of this study with the 8-view setting. In this study, we additionally focused on the \mathcal{L}_{KL} loss. In particular, we both tested the same implementation of the KL-Divergence loss used in

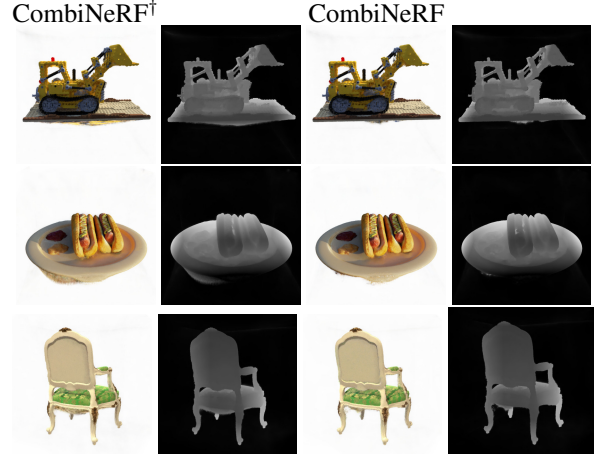


Figure 7. Qualitative result on the ablation study of NeRF-Synthetic with 8-view setting. We compare our CombiNeRF with CombiNeRF[†] (our method with KL-Divergence as in InfoNeRF).

[10] (implemented in an intermediate version of CombiNeRF called CombiNeRF[†]) and our modified version (as described in Sec. 3.2) used in the final CombiNeRF.

In Fig. 7 we show the qualitative results of both the two implementations. While quantitative results remain almost the same (see Tab. 4), CombiNeRF is able to remove most of the artifacts generated below the objects, thus indicating an improved reconstruction of the final geometries.

5. Conclusion

In this work, we presented CombiNeRF, a combination of regularization techniques in a unified framework. We regularize neighboring rays distributions, we also take into account the single ray distribution and a smoothness term is adopted to regularize near geometries. Besides these additional loss functions, we also consider Lipschitz regularization and an encoding mask to regularize high-frequency components.

CombiNeRF shows cutting-edge and consistent results in the quality of the reconstructions and it outperforms the SOTA methods in LLFF and NeRF-Synthetic datasets with few-shot settings. The increasing performance of the combination of multiple regularization techniques is validated through an ablation study that highlights the contribution of the CombiNeRF components and the differences with the previous implementations.

Additional techniques like Ray Entropy Loss [10], Adversarial Perturbation [5] and Space Annealing [23] have been taken into account but not considered in CombiNeRF for their apparently poor contribution. Additional details are given in the supplementary material but more experiments are required as future works on these and other techniques to achieve a more general purpose NeRF framework.

References

- [1] Benjamin Attal, Jia-Bin Huang, Michael Zollhöfer, Johannes Kopf, and Changil Kim. Learning neural light fields with ray-space embedding networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [2] Jiayang Bai, Letian Huang, Wen Gong, Jie Guo, and Yanwen Guo. Self-nerf: A self-training pipeline for few-shot neural radiance fields, 2023. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5470–5479, 2022. 1, 2, 4
- [4] Eric R. Chan, Koki Nagano, Matthew A. Chan, Alexander W. Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. GeNVS: Generative novel view synthesis with 3D-aware diffusion models. In *arXiv*, 2023. 2
- [5] Tianlong Chen, Peihao Wang, Zhiwen Fan, and Zhangyang Wang. Aug-nerf: Training stronger neural radiance fields with triple-level physically-grounded augmentations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 8, 1
- [6] Guangcong, Zhaoxi Chen, Chen Change Loy, and Ziwei Liu. Sparsenerf: Distilling depth ranking for few-shot novel view synthesis. *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 2
- [7] Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. Nerfren: Neural radiance fields with reflections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18409–18418, 2022. 2
- [8] Shoukang Hu, Kaichen Zhou, Kaiyu Li, Longhui Yu, Lanqing Hong, Tianyang Hu, Zhenguo Li, Gim Hee Lee, and Ziwei Liu. Consistentnerf: Enhancing neural radiance fields with 3d consistency for sparse view synthesis, 2023. 2
- [9] 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 (ICCV)*, pages 5885–5894, 2021. 5, 2
- [10] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *CVPR*, 2022. 1, 2, 3, 8
- [11] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 4
- [12] Stefan Lionar, Lukas Schmid, Cesar Cadena, Roland Siegwart, and Andrei Cramariuc. Neuralblox: Real-time neural representation fusion for robust volumetric mapping. *2021 International Conference on 3D Vision (3DV)*, 2021. 2
- [13] Hsueh-Ti Derek Liu, Francis Williams, Alec Jacobson, Sanja Fidler, and Or Litany. Learning smooth neural functions via lipschitz regularization. *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings*, 2022. 1, 2, 5
- [14] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, 2019. 2
- [15] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 5
- [16] 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. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2, 3, 4, 5
- [17] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P. Srinivasan, and Jonathan T. Barron. NeRF in the dark: High dynamic range view synthesis from noisy raw images. *CVPR*, 2022. 2
- [18] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. 1, 2, 3, 5
- [19] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty 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. 2
- [20] Phong Nguyen-Ha, Animesh Karnewar, Lam Huynh, Esa Rahtu, and Janne Heikkilä. Rgb-d-net: Predicting color and depth images for novel views synthesis. In *Proceedings of the International Conference on 3D Vision*, 2021. 2
- [21] Michael Niemeyer and Andreas Geiger. Campari: Camera-aware decomposed generative neural radiance fields. In *International Conference on 3D Vision (3DV)*, 2021. 2
- [22] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [23] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 4, 5, 8
- [24] M. Píala and R. Clark. Terminerf: Ray termination prediction for efficient neural rendering. In *2021 International Conference on 3D Vision (3DV)*, pages 1106–1114, Los Alamitos, CA, USA, 2021. IEEE Computer Society. 2
- [25] Matteo Poggi, Pierluigi Zama Ramirez, Fabio Tosi, Samuele Salti, Luigi Di Stefano, and Stefano Mattoccia. Cross-spectral neural radiance fields. In *Proceedings of the International Conference on 3D Vision*, 2022. 3DV. 2
- [26] 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*, 2020. 2
- [27] Radu Alexandru Rosu and Sven Behnke. Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 5
- [28] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- [29] Kevin Scaman and Aladin Virmaux. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *Neural Information Processing Systems*, 2018. 1
- [30] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [31] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2019. 2
- [32] Liangchen Song, Zhong Li, Xuan Gong, Lele Chen, Zhang Chen, Yi Xu, and Junsong Yuan. Harnessing low-frequency neural fields for few-shot view synthesis, 2023. 2
- [33] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. P. Mildenhall, P. Srinivasan, J. T. Barron, and H. Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8238–8248, Los Alamitos, CA, USA, 2022. IEEE Computer Society. 2
- [34] Jiaxiang Tang. Torch-ngp: A pytorch implementation of instant-ngp, 2022. 1, 5
- [35] 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. *CVPR*, 2022. 2
- [36] Jingwen Wang, Tymoteusz Bleja, and Lourdes Agapito. Go-surf: Neural feature grid optimization for fast, high-fidelity rgb-d surface reconstruction. *2022 International Conference on 3D Vision (3DV)*, 2022. 2
- [37] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *CoRR*, abs/2106.10689, 2021. 2
- [38] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and Wenping Wang. F2-nerf: Fast neural radiance field training with free camera trajectories. *CVPR*, 2023. 2
- [39] Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. Hf-neus: Improved surface reconstruction using high-frequency details. *Advances in Neural Information Processing Systems*, 35:1966–1978, 2022. 2
- [40] 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. 5
- [41] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yu-Xiong Wang, and David Forsyth. Diver: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [42] Jamie Wynn and Daniyar Turmukhambetov. DiffusioNeRF: Regularizing Neural Radiance Fields with Denoising Diffusion Models. In *CVPR*, 2023. 1, 2, 4, 5
- [43] Christopher Xie, Keunhong Park, Ricardo Martin-Brualla, and Matthew Brown. Fig-nerf: Figure-ground neural radiance fields for 3d object category modelling. In *International Conference on 3D Vision (3DV)*, 2021. 2
- [44] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2023. 4, 5
- [45] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. 2
- [46] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [47] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. PhysSG: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [48] Kai Zhang, Fujun Luan, Zhengqi Li, and Noah Snavely. Iron: Inverse rendering by optimizing neural sdfs and materials from photometric images. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2
- [49] 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. 5
- [50] Fuqiang Zhao, Yuheng Jiang, Kaixin Yao, Jiakai Zhang, Liao Wang, Haizhao Dai, Yuhui Zhong, Yingliang Zhang, Minye Wu, Lan Xu, and Jingyi Yu. Human performance modeling and rendering via neural animated mesh. *ACM Trans. Graph.*, 41(6), 2022. 2
- [51] Yihao Zhi, Shenhan Qian, Xinhao Yan, and Shenghua Gao. Dual-space nerf: Learning animatable avatars and scene lighting in separate spaces. In *International Conference on 3D Vision (3DV)*, 2022. 2

CombiNeRF: A Combination of Regularization Techniques for Few-Shot Neural Radiance Field View Synthesis

Supplementary Material

6. Additional Implementations

In the following sections, we describe three additional implementations considered during the empirical selection of the best methods used in CombiNeRF. The following implementations were not included in CombiNeRF for their apparently poor contribution with respect to Vanilla NeRF [18]. Due to a lack of time for further testing, we were only able to partially evaluate these methods, which will eventually be considered in future CombiNeRF extensions. As done for the other considered techniques, we initially tested their performance on the LLFF dataset with the 3-view setting, comparing their results against Vanilla NeRF. For completeness, we provide below a brief summary of these techniques.

6.1. Regularization Losses

Ray Entropy Loss. The aim is to introduce a sparsity constraint to focus only on the scene of interest. This constraint can be achieved by minimizing the entropy of each sampled ray density function, as done in [10], in order to force the rays to have a few high peaks, i.e. the peaks should only represent the hit surfaces.

To minimize the entropy of a ray \mathbf{r} , we start defining the normalized ray density $q(\mathbf{r})$ as:

$$q(\mathbf{r}_i) = \frac{\alpha_i}{\sum_j \alpha_j} = \frac{1 - \exp(-\sigma_i \delta_i)}{\sum_j 1 - \exp(-\sigma_j \delta_j)}, \quad (15)$$

where \mathbf{r}_i ($i = 1, \dots, N$) is a sampled point in the ray, σ_i is the observed density at \mathbf{r}_i , δ_i is a sampling interval around \mathbf{r}_i and α_i is the opacity at \mathbf{r}_i .

Starting from the Shannon Entropy, we can define the entropy of a discrete ray density function:

$$H(\mathbf{r}) = - \sum_{i=1}^N q(\mathbf{r}_i) \log q(\mathbf{r}_i), \quad (16)$$

where the calculation of $q(\mathbf{r}_i)$, which includes σ_i and δ_i , is already given thanks to volume rendering computation.

A mask M is applied to filter the rays in the entropy defined in Eq. (16) with the aim of discarding too low-density rays not intersecting or hitting any object in the scene:

$$M(\mathbf{r}) = \begin{cases} 1 & \text{if } Q(\mathbf{r}) > \epsilon \\ 0 & \text{otherwise} \end{cases}, \quad (17)$$

where

$$Q(\mathbf{r}) = \sum_{i=1}^N 1 - \exp(-\sigma_i \delta_i), \quad (18)$$

is the cumulative ray density.

Because unobserved viewpoints can help generalization, for the entropy loss computation both rays from training views and from unseen images can be considered, whose sets are denoted respectively by R_t and R_u . The final entropy loss is defined as:

$$\mathcal{L}_{entr} = \lambda_{entr} \cdot \frac{1}{|R_t| + |R_u|} \sum_{\mathbf{r} \in R_t \cup R_u} M(\mathbf{r}) \odot H(\mathbf{r}), \quad (19)$$

where \odot denotes the element-wise multiplication and λ_{entr} defines the contribution of this loss.

Adversarial Perturbation. Previous works demonstrate how neural networks benefit from random or learned data augmentation. Generalization can be achieved by injecting noise during the training phase, as done in [5]. Here, noise is introduced in the form of worst-case perturbations applied to the input coordinates, features of the network and the pre-rendering output of the network. The worst-case perturbation can be formulated as a min-max problem defined as follows:

$$\min_{\Theta} \max_{\delta} \|C^\dagger(r|\Theta, \delta) - C^*(r)\|_2^2, \quad (20)$$

where

$$\delta = (\delta_p, \delta_f, \delta_r) \in \mathcal{S}_p \times \mathcal{S}_f \times \mathcal{S}_r, \quad (21)$$

and $\delta_p, \delta_f, \delta_r$ are the learned perturbation, $\mathcal{S}_p, \mathcal{S}_f, \mathcal{S}_r$ are the corresponding search range and $C^\dagger(r|\Theta, \delta)$ is the perturbed rendered color.

The input coordinates perturbation $\delta_p = (\delta_t, \delta_{xyz}, \delta_\theta)$ consists of along-ray perturbation $\delta_t \in \mathbb{R}^3$, point position perturbation $\delta_{xyz} \in \mathbb{R}^3$ and view-direction perturbation $\delta_\theta \in \mathbb{R}^3$ (hence $\mathcal{S}_p \subseteq \mathbb{R}^6$). The feature perturbation δ_f consists of perturbing a D -dimensional features vector, hence $\mathcal{S}_f \subseteq \mathbb{R}^D$. Finally, a pre-rendering output perturbation is applied to RGB color and density values ($\mathcal{S}_r \subseteq \mathbb{R}^4$).

In order to estimate the worst-case perturbation of the injected noise, [5] proposes a multi-step Projected Gradient Descent (PGD) approach. Given l_{\inf} as the norm ball defining all search spaces $\mathcal{S}_p, \mathcal{S}_f, \mathcal{S}_r$ and the radius ϵ as the maximum magnitude of the perturbation, a PGD step is defined as:

$$\delta^{(t+1)} = \prod_{\|\delta\|_\infty \leq \epsilon} \left[\delta^{(t)} + \alpha \cdot \text{sgn}(\|C^\dagger(r|\Theta, \delta) - C^*(r)\|_2^2) \right], \quad (22)$$

LLFF 3-views	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
Vanilla NeRF	17.71	0.544	0.303	0.155
\mathcal{L}_{adv}	17.75	0.541	0.300	0.155
\mathcal{L}_{entr}	17.34	0.526	0.312	0.161
Anneal	18.04	0.562	0.284	0.147
CombiNeRF	20.37	0.686	0.191	0.101

Table 5. Quantitative results comparison among the other three regularization techniques on LLFF dataset with 3-view setting.

where α is the step size, $\Pi[\cdot]$ is a projection operator and $\text{sgn}(\cdot)$ is the sign function.

Following this procedure, the adversarial perturbation loss is obtained as:

$$\mathcal{L}_{adv} = \lambda_{adv} \cdot \sum_{r \in \mathcal{R}} \|C^\dagger(r|\Theta, \delta) - C^*(r)\|_2^2, \quad (23)$$

where λ_{adv} defines the contribution of this loss.

6.2. Space Annealing

When few images are provided, NeRF could converge to undesired solutions. In this case, high-density values are concentrated close to ray origins. Annealing the sampled scene space over the early iterations of the training [23] can be a solution to avoid this specific problem: the idea is to restrict the sampling space to an initial smaller region where the scene is centered.

Given the camera’s near and far planes t_n and t_f , let t_m be a point between them, likely their midpoint, and define the new near and far plane per-iteration as:

$$\begin{aligned} t_n(i) &= t_m + (t_n - t_m)\eta(i) \\ t_f(i) &= t_m + (t_f - t_m)\eta(i) \\ \eta(i) &= \min(\max(i/N_t, p_s), 1), \end{aligned} \quad (24)$$

where i is the current training iteration, the parameter N_t represents the maximum number of iterations in which to perform the space annealing and p_s indicates the starting range. These parameters should be tuned according to the particular scene under consideration.

Annealing the sampled space, especially during the first iteration of training, may allow NeRF to focus on the right region of interest, avoiding degenerate solutions.

6.3. Results

We evaluated the three implementations presented above on the LLFF dataset with a 3-view setting. We compared the obtained results with the results obtained by Vanilla NeRF to see if a specific technique should be considered for further experiments or should be revised or better fine-tuned.

In Tab. 5 we show the quantitative results of this study. CombiNeRF remains the best-performing method, outperforming all the other methods by a large margin. \mathcal{L}_{entr} is not able to improve NeRF generalization and, on the contrary, degrades its performance. \mathcal{L}_{adv} slightly improves some metrics but it doesn’t affect the overall score. Moreover, the time complexity required to perform PDG for worst-case perturbation considerably increases the total computational time. Instead, the Space Annealing technique outperforms Vanilla NeRF in all metrics. This result would lead us to focus on the latter as a promising method to be tested and eventually integrated into CombiNeRF. However, we argue that the contribution provided by Space Annealing could not offer additional benefits when paired with the distortion loss of Eq. (8), as both methods try to remove high-density values on points close to the camera.

In Fig. 8 we show qualitative results of Vanilla NeRF against the other three implementations. In Space Annealing results, the depths are smoother and the reconstruction quality increases as visible in the white wall and the light on top of "T-rex", in the leaves’ details on "Fern" and in the overall object reconstruction of "Fortress". With \mathcal{L}_{adv} is very difficult to see improvements in the resulting renderings, which reflects the quantitative results seen in Tab. 5. In \mathcal{L}_{entr} results, both depth and RGB images are less accurate, containing more artifacts (e.g., in "Fortress") and noise, particularly visible in the white wall of "T-rex".

7. Training Details

For an easier reproduction of the results, we provide the training procedure adopted for the LLFF dataset, which follows [23], and for the NeRF-Synthetic dataset, which follows [9].

7.1. LLFF

In LLFF we use all 8 scenarios ("Fern", "Room", "T-rex", "Flower", "Leaves", "Horns", "Orchids" and "Fortress"). Each view of each scene is 378×504 , thus $8 \times$ downsampled with respect to the original resolution. Test views are taken every 8th view and input images are evenly sampled among the remaining ones.

7.2. NeRF-Synthetic

In Nerf-Synthetic we use all 8 scenarios ("Lego", "Hot-dog", "Mic", "Drums", "Materials", "Ship", "Chair" and "Ficus"). Each view of each scene is 400×400 , thus $2 \times$ downsampled with respect to the original resolution. We take 25 test views evenly sampled from the original test set, while training views are chosen according to the following IDs: 2, 16, 26, 55, 73, 75, 86, 93.

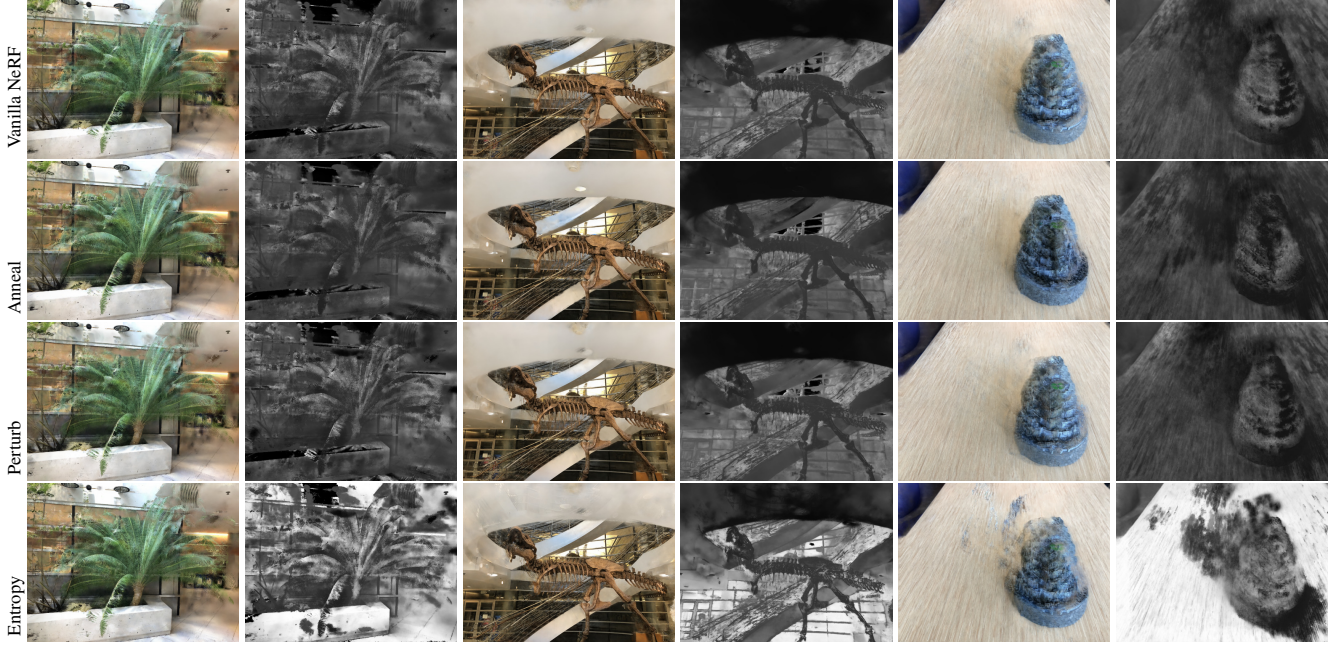


Figure 8. Comparison of Vanilla NeRF against Space Annealing, Adversarial Perturbation and Entropy loss on LLFF with 3-view setting.

8. Additional Results

In this section we show additional quantitative and qualitative results of CombiNeRF on LLFF with 3/6/9-views setting and NeRF-Synthetic dataset with 8-view setting. Evaluation results are taken from the respective paper, when reported.

8.1. Qualitative Results

In Fig. 9, Fig. 10 and Fig. 11 we show additional qualitative results on all scenarios of the LLFF dataset under 3-view, 6-view, 9-view settings, respectively. Under the 3-view setting, we can see the effectiveness of CombiNeRF against Vanilla NeRF in both RGB images and depth images. When the number of input images increases, the performance gap between Vanilla NeRF and CombiNeRF decreases, but the depth quality still remains definitely better in the latter.

Fig. 12 shows additional qualitative results on all scenarios of the NeRF-Synthetic dataset under the 8-view setting. CombiNeRF is able to better reconstruct the scenarios, and this is particularly visible in the geometry of "Ship" and in the color of "Ficus". Finer details are less visible in Vanilla NeRF, as we can see in "Mic" and "Chair", while a lot of floaters and noise are removed in "Drums" and "Materials" with CombiNeRF.

8.2. Quantitative Results

In Tab. 6, Tab. 7 and Tab. 8 we show per-scene quantitative results on PSNR, SSIM, and LPIPS metrics, respectively. We can observe that in "Lego", "Chair", "Mic",

"Hotdog" and "Ficus" CombiNeRF gets the better results in all the metrics. Instead in "Drums", "Materials" and "Ship", which are the most challenging scenarios, CombiNeRF outperforms the other methods in the overall metric scores.

In Tab. 9, Tab. 10 and Tab. 11 we show quantitative results on PSNR, SSIM, and LPIPS metrics respectively, including 3/6/9-view settings. CombiNeRF clearly outperforms Vanilla NeRF in all the scenarios.

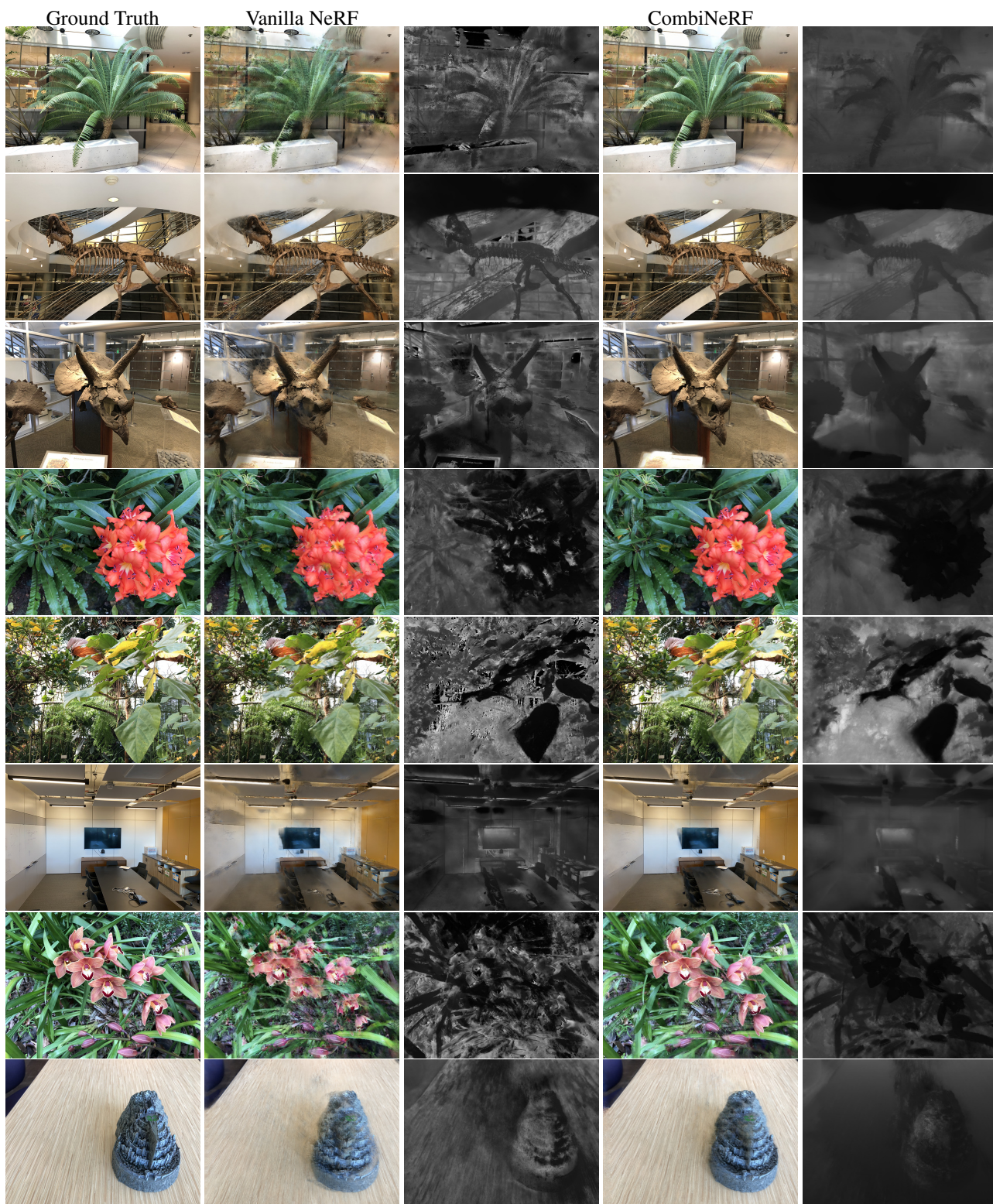


Figure 9. Additional qualitative results on LLFF dataset with 3-view setting.

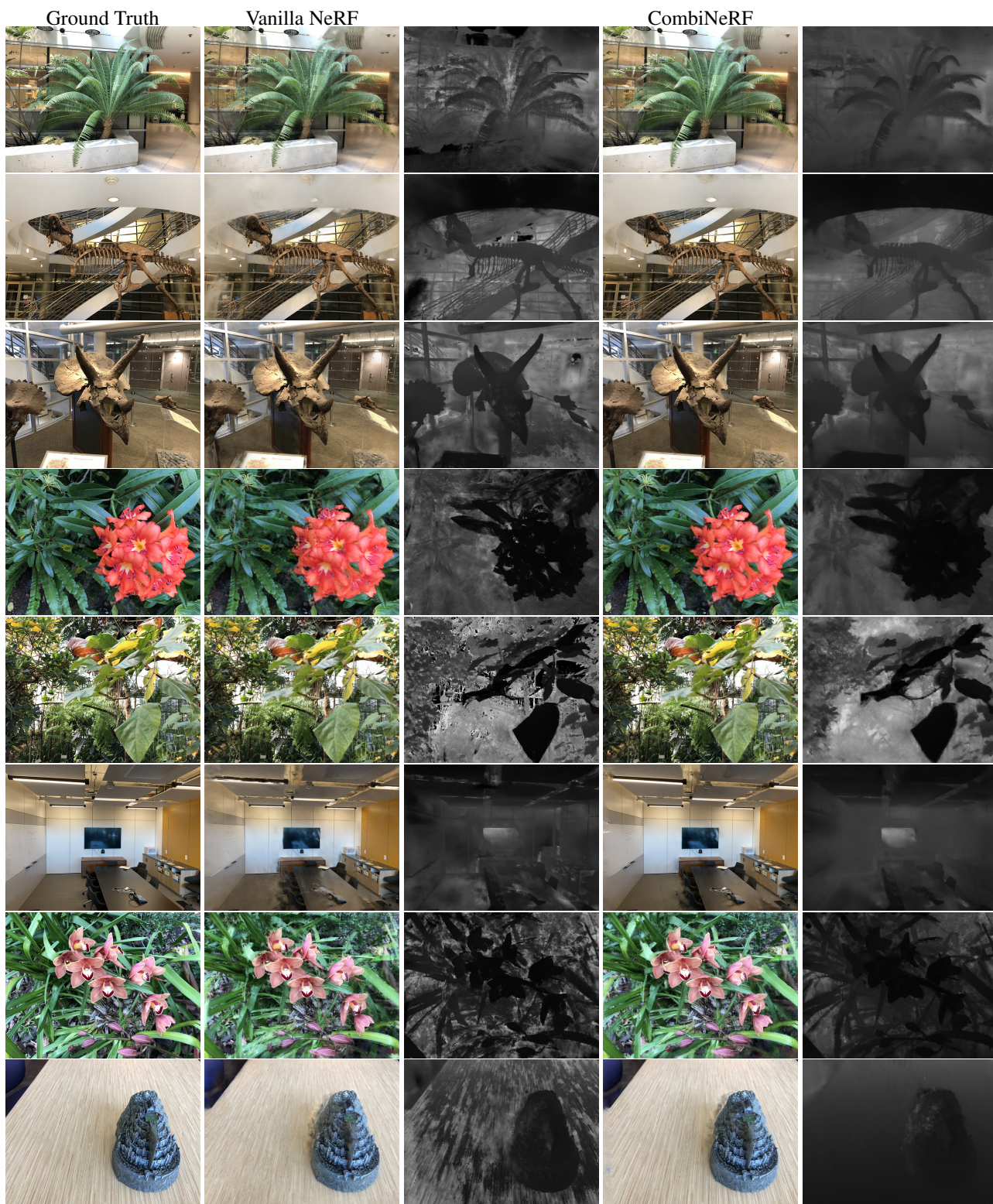


Figure 10. Additional qualitative results on LLFF dataset with 6-view setting.

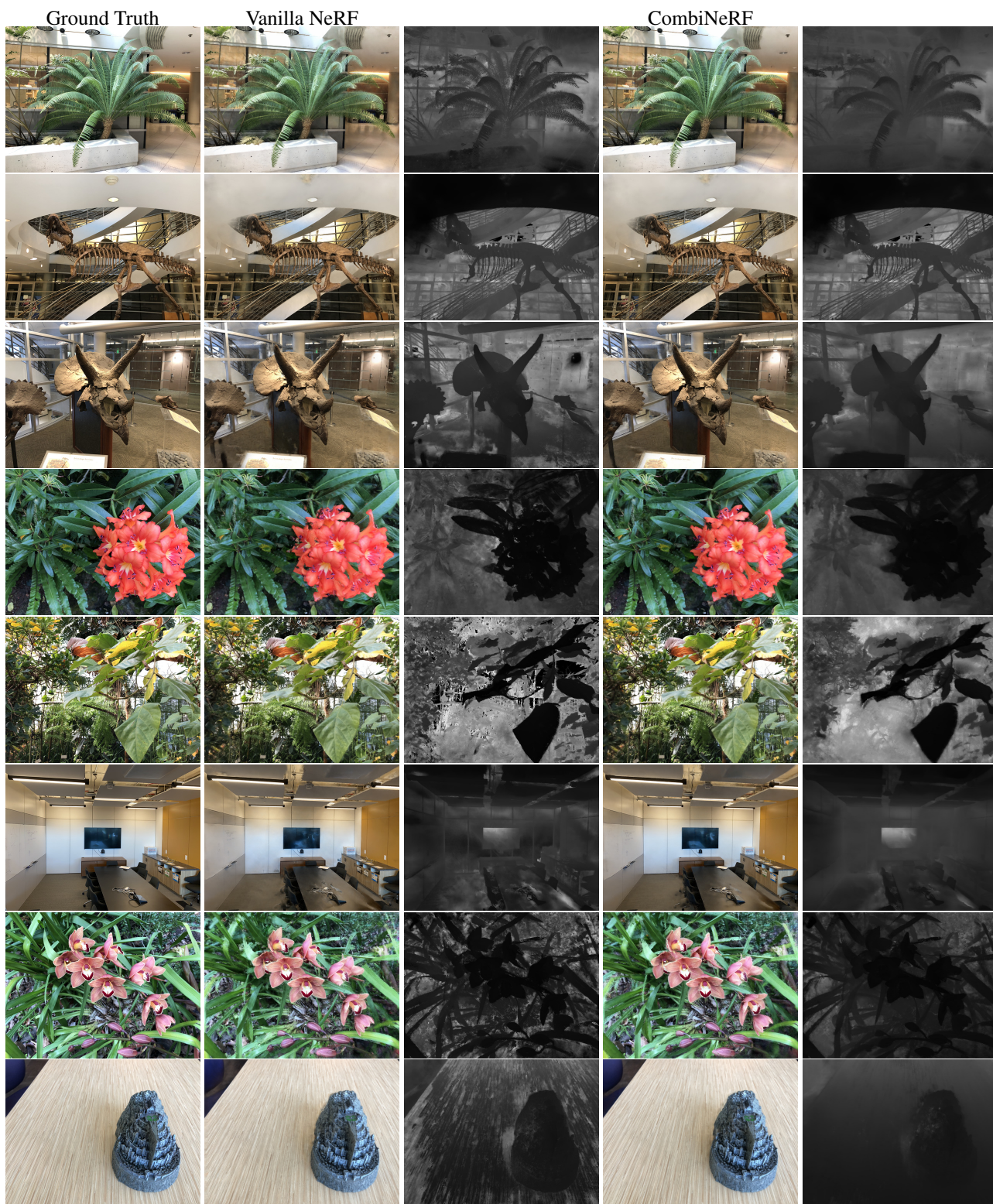


Figure 11. Additional qualitative results on LLFF dataset with 9-view setting.

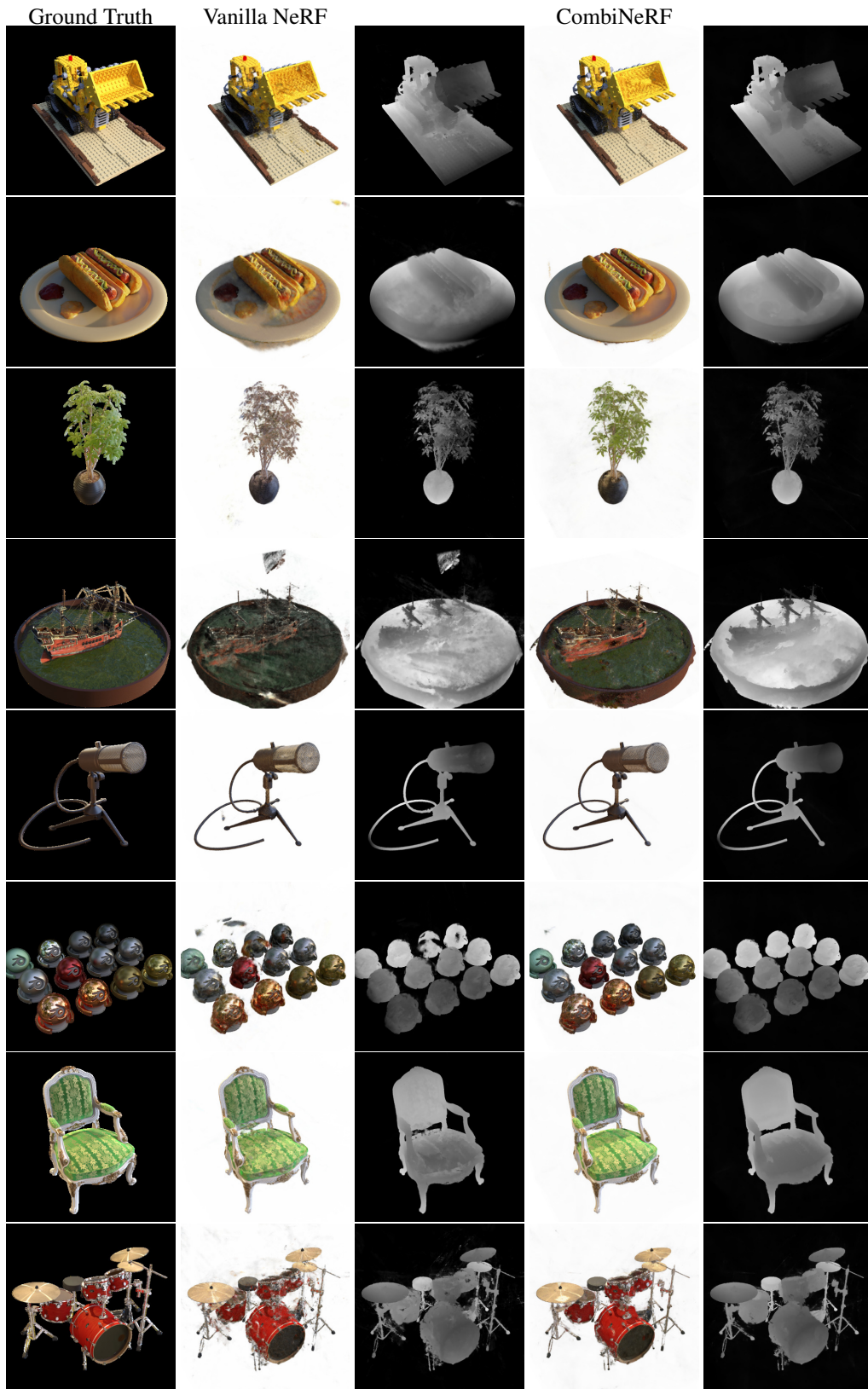


Figure 12. Additional qualitative results on NeRF-Synthetic dataset with 8-view setting.

NeRF-Synthetic 8-views	Drums	Material	Ship	Ficus	Lego	Chair	Mic	Hotdog
Vanilla NeRF	18.07	19.36	18.69	22.93	22.49	25.25	27.34	24.55
DietNeRF, \mathcal{L}_{MSE} ft	20.029	21.621	22.536	20.940	24.311	25.595	26.794	26.626
CombiNeRF	20.165	20.73	21.392	23.493	24.958	27.862	28.172	28.383

Table 6. Per-scene quantitative results on NeRF-Synthetic dataset w.r.t. **PSNR** \uparrow metric.

NeRF-Synthetic 8-views	Drums	Material	Ship	Ficus	Lego	Chair	Mic	Hotdog
Vanilla NeRF	0.767	0.804	0.69	0.897	0.851	0.901	0.949	0.901
DietNeRF, \mathcal{L}_{MSE} ft	0.845	0.851	0.757	0.874	0.875	0.912	0.950	0.924
CombiNeRF	0.838	0.838	0.754	0.91	0.885	0.933	0.957	0.945

Table 7. Per-scene quantitative results on NeRF-Synthetic dataset w.r.t. **SSIM** \uparrow metric.

NeRF-Synthetic 8-views	Drums	Material	Ship	Ficus	Lego	Chair	Mic	Hotdog
Vanilla NeRF	0.222	0.185	0.269	0.107	0.105	0.076	0.054	0.134
DietNeRF, \mathcal{L}_{MSE} ft	0.117	0.095	0.193	0.094	0.096	0.077	0.043	0.067
CombiNeRF	0.14	0.114	0.151	0.078	0.072	0.049	0.035	0.066

Table 8. Per-scene quantitative results on NeRF-Synthetic dataset w.r.t. **LPIPS** \downarrow metric.

(a) 3 input views.								
LLFF 3-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	12.47	17.72	15.65	18.54	20.54	19.52	17.79	19.47
CombiNeRF	16.04	18.46	18.84	21.01	21.21	21.47	22.21	23.75
(b) 6 input views.								
LLFF 6-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	16.74	20.1	22.32	22.25	21.66	25.58	23.18	24.42
CombiNeRF	18.11	20.46	23.26	24.52	23.91	28.72	25.07	27.9
(c) 9 input views.								
LLFF 9-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	18.45	21.07	24.59	24.88	25.27	27.11	25.54	26.76
CombiNeRF	19.1	21.22	25.15	26.13	26.35	28.57	26.17	28.51

Table 9. Per-scene quantitative results on LLFF dataset w.r.t. **PSNR** \uparrow metric.

(a) 3 input views.								
LLFF 3-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	0.199	0.61	0.523	0.548	0.719	0.755	0.509	0.486
CombiNeRF	0.461	0.666	0.692	0.668	0.773	0.821	0.713	0.692
(b) 6 input views.								
LLFF 6-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	0.499	0.738	0.797	0.725	0.816	0.894	0.744	0.688
CombiNeRF	0.584	0.753	0.823	0.815	0.863	0.921	0.811	0.873
(c) 9 input views.								
LLFF 9-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	0.605	0.777	0.86	0.82	0.885	0.912	0.827	0.805
CombiNeRF	0.647	0.782	0.876	0.857	0.908	0.93	0.85	0.877

Table 10. Per-scene quantitative results on LLFF dataset w.r.t. **SSIM** \uparrow metric.

(a) 3 input views.								
LLFF 3-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	0.435	0.181	0.35	0.3	0.195	0.269	0.402	0.292
CombiNeRF	0.25	0.155	0.212	0.226	0.15	0.184	0.194	0.157
(b) 6 input views.								
LLFF 6-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	0.23	0.124	0.137	0.156	0.133	0.112	0.166	0.137
CombiNeRF	0.178	0.12	0.122	0.095	0.094	0.08	0.105	0.054
(c) 9 input views.								
LLFF 9-view	Orchids	Leaves	Horns	Flower	T-rex	Room	Fern	Fortress
Vanilla NeRF	0.171	0.106	0.091	0.093	0.078	0.083	0.101	0.084
CombiNeRF	0.151	0.112	0.078	0.071	0.06	0.066	0.082	0.05

Table 11. Per-scene quantitative results on LLFF dataset w.r.t. **LPIPS** \downarrow metric.