# 4K-NeRF: High Fidelity Neural Radiance Fields at Ultra High Resolutions

Zhongshu Wang, Lingzhi Li, Zhen Shen, Li Shen, Liefeng Bo

Alibaba Group

Beijing, China

{zhongshu.wzs, llz273714, zackary.sz, jinyan.sl, liefeng.bo}@alibaba-inc.com

## Abstract

*In this paper, we present a novel and effective framework, named 4K-NeRF, to pursue high fidelity view synthesis on the challenging scenarios of ultra high resolutions, building on the methodology of neural radiance fields (NeRF). The rendering procedure of NeRF-based methods typically relies on a pixel wise manner in which rays (or pixels) are treated independently on both training and inference phases, limiting its representational ability on describing subtle details especially when lifting to a extremely high resolution. We address the issue by better exploring ray correlation for enhancing high-frequency details benefiting from the use of geometry-aware local context. Particularly, we use the view-consistent encoder to model geometric information effectively in a lower resolution space and recover fine details through the view-consistent decoder, conditioned on ray features and depths estimated by the encoder. Joint training with patch-based sampling further facilitates our method incorporating the supervision from perception oriented regularization beyond pixel wise loss. Quantitative and qualitative comparisons with modern NeRF methods demonstrate that our method can significantly boost rendering quality for retaining high-frequency details, achieving the state-of-the-art visual quality on 4K ultra-high-resolution scenario. Code Available at* https://github.com/frozoul/4K-NeRF

## 1. Introduction

Ultra-High-Resolution has growing popular as a standard for recording and displaying images and videos, even supported in modern mobile devices. A scene captured in ultra high resolution format typically presents content incredible details compared to using a relatively lower resolution (e.g, 1K high-definition format) in which the information of a pixel is enlarged by a small patch in the extremely high resolution images. Developing techniques for handling such high-frequency details poses challenges for a wide range of tasks in image processing and computer vision. In this paper,



Figure 1. **Qualitative comparison between our method (left) and DVGO (right).** The picture is better displayed in a high resolution display for visualization.

we focus on the novel view synthesis task and investigate the potential of realizing high fidelity view synthesis rich in subtle details at ultra high resolution.

Novel view synthesis aims to produce free-view photo-realistic synthesis given a sparse set of images captured for a scene from multiple viewpoints. Recently, Neural Radiance Fields [26] offer a new methodology for modelling and rendering 3D scenes by virtue of deep neural networks and have demonstrated remarkable success on improving visual quality compared to traditional view interpolation methods [37, 37, 45]. Particularly, a mapping function, instantiated as a deep multilayer perceptron (MLP), is optimized to associate each 3D location given a viewing direction to its corresponding radiance color and volume density, while realizing view-dependent effect requires querying the large network hundreds of times for the ray casting through each pixel. Several following approaches are proposed to improve the method either from the respect of reducing aliasing artifacts on multiple scales [1] or improving training and inference efficiency benefiting from the use of discretized

structures [6, 36, 46]. All these methods follow the pixel-wise mechanism despite varying architectures, i.e., rays (or pixels) are regarded individually during training and inference phase. They are typically developed on training views up to 1K resolution. When applying the approaches on ultra-high-resolution scenarios, they would struggle with objectionable blurring artifacts (as shown in Fig. 3) due to insufficient representational ability for capturing fine details.

In this paper, we introduce a novel framework, named 4K-NeRF, building on the methodology of NeRF-based volume rendering to realize high fidelity view synthesis at 4K ultra high resolution. We take the inspiration from the success of convolutional neural networks on traditional super resolution [11], which serves to resolve a lower resolution observation to a higher resolution with rich details by virtue of local priors learned between neighbouring pixels. We expect to boost the representational power of NeRF-based methods by better exploring local correlations between rays.

Specially, the framework is comprised of two components, a view-consistent encoder and a view-consistent decoder, as shown in Fig. 6. The encoder is exploited to encode geometric properties of a scene effectively in a lower resolution space, forming intermediate ray features and geometry information (i.e., estimated depth) feeding into the decoder. The decoder is capable of recovering high-frequency details by integrating geometry-aware local patterns learned via depth-modulated convolutions in the higher resolution (full-scale) observations. We further introduce a patch-based ray sampling strategy replacing the random sampling in NeRFs, allowing the encoder and the decoder trained jointly with the perception oriented losses (i.e., adversarial loss and perceptual loss) complementing to the conventional pixel-wise MSE loss. More importantly, such jointly training assists geometric modelling in the encoder coordinated with the learning of local context in the decoder, realizing view-consistent enhancement on fine details. Empirical comparisons and ablation studies on challenging scenarios with 4K ultra high resolution demonstrate the effectiveness of the proposed framework both quantitatively and qualitatively. We further validate the generalization of the method on different base architectures as well as improving visual quality on the standard 1K resolution setting.

## 2. Related work

**Novel View Synthesis.** Novel view synthesis is a well-studied task. Before the rise of implicit representation, many effort have been done by the research community, the mainstream methods included mesh-based [16, 32], 3D geometry-based [8, 14, 35],point cloud based [43], multiplane images-based [4, 12, 37], etc. However, because they still try to reconstruct a kind of geometry representation, e.g., a mesh or multi-plane images, these methods still have difficulty handling some extremely difficult scenes that have complex geometry and large camera gap.

**Neural Radiance Fields.** The emerging neural radiance field marks the rise of a new paradigm for novel-view synthesis. NeRF directly learns a continuous mapping from 3D coordinates and view directions to view-dependent color and volume density and acquires pixel color through the volume rendering technique. It can synthesize photo-realistic novel views with out direct 3D supervision on both real forward-facing scenes and synthetic bounded scenes and shows robust scalability. [13, 15, 17, 20, 29, 47] accelerated rendering speed from originally multi-seconds to milliseconds level. [6, 23, 28, 36, 46] introduced volumetric radiance fields to successfully reduce training cost with orders of magnitude. Some methods focus on the aspect of improving the rendering quality of NeRF [50]. [1, 2] introduced mipmap for achieving anti-aliasing, and [38, 44] improved NeRF's ability on modelling high reflectance surface, .The approaches proposed in [10, 30, 48] incorporate depth prior or 2D image prior to reducing cloudy artifacts caused by insufficient training views. [39] adopted a super resolution module directly to improve visual quality relying on high resolution guidance during inference while the method might still encounter the issue of inconsistent view rendering.

To the best of our knowledge, our framework is the first to successfully extend NeRF-based paradigm to 4K resolution, proving high-fidelity viewing experience with crystal-clear and high-frequency details.

**High-Resolution Synthesis.** Image super-resolution is a problem that has been studied for a long time, aiming to recover a high-resolution image from a single low-resolution image. The classical super-resolution methods assume that the image degradation process only contains downsampling and noise [33]. Under this setting, there have been many super-resolution methods that can achieve excellent results [7, 11, 19, 22, 42, 52]. However, in real-world scenes, image degradation is often affected by many factors, such as noisy and blurry artifacts, compression distortion, and the combinations in a different order. Classical super-resolution method is sometimes less effective in practice. In order to solve the problem, some real-world super-resolution methods have been proposed recently [18, 21, 40, 49]. These methods achieved satisfactory results in real-world images. All of these super-resolution methods perform on resolving 2D single image.

## 3. Method

We first go over the methodology of NeRF-based volumetric rendering and discuss the limitation on modeling and rendering the scenes with extremely high resolution. We then present our NeRF-4K framework in detail and introduce the training strategy with loss functions in the next section.
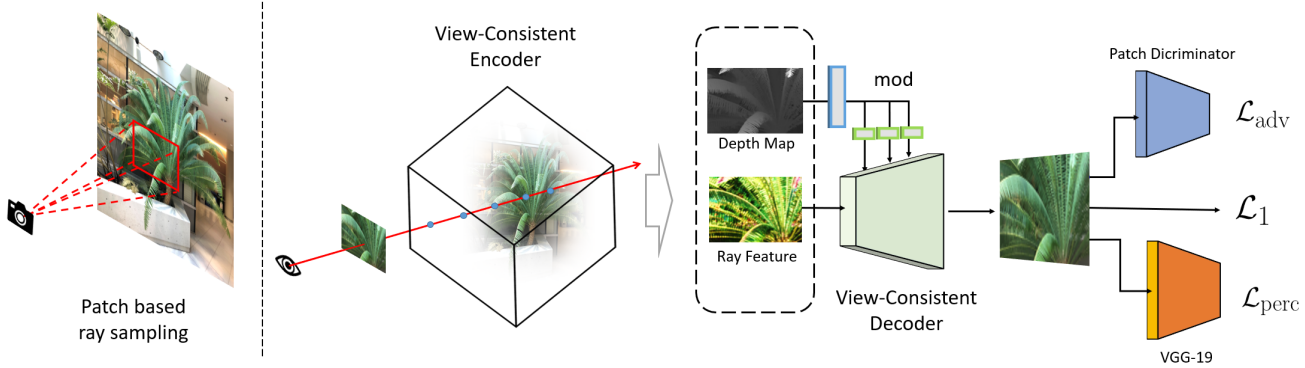
Figure 2. **The overall pipeline of 4K-NeRF.** Using patch-based ray sampling, we jointly train the VC-Encoder for encoding 3D geometric information in a lower resolution space and the VC-Decoder for realizing high quality rendering with view-consistent enhancement on high-frequency details.

## 3.1. Volumetric Rendering

NeRF realizes photo-realistic view synthesis by learning a continuous mapping function to estimate the color $\mathbf{c} \in \mathbb{R}^3$ and the volume density $\sigma \in \mathbb{R}$ of a 3D point position $\mathbf{Fx} \in \mathbb{R}^3$ and a viewing direction $\mathbf{d} \in \mathbb{R}^3$, i.e., $\Phi : (\mathbf{x}, \mathbf{d}) \mapsto (\mathbf{c}, \sigma)$. To render an image given camera pose, the expected color $\widehat{\mathbf{C}}(\mathbf{r})$ of a camera ray $\mathbf{r} = \mathbf{o} + t\mathbf{d}$ from the camera center $\mathbf{o}$ through the pixel is estimated with the the numerical quadrature discussed by Max [24], by sampling a set of points along the ray and integrating their colors to approximate a volumetric rendering integral,

$$\widehat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^{N} T_i \cdot \alpha_i \cdot \mathbf{c}_i, \qquad (1)$$

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j), \qquad (2)$$

where $\alpha_i$ denotes the ray termination probability at the point $i$, $\delta_i = t_{i+1} - t_i$ represents the distance between two adjacent points, and $T_i$ indicates the accumulated transmittance when reaching $i$. The mapping function $\Phi$ is instantiated as a multilayer perceptron (MLP). Given the training set of images with known poses, the model is trained by minimizing the mean squared errors (MSE) between the predicted pixel colors and the ground-truth colors,

$$\mathcal{L}_{\text{MSE}} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \left\| \widehat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2, \qquad (3)$$

where $\mathcal{R}$ denotes the ray set randomly sampled in each minibatch. The optimization of each point is according to its projection through the rays of different viewpoints. Some variants are proposed to integrate the learning with explicit structures [23, 29, 36, 46] instead of single large neural networks. By confining density prediction to the use of po-

sitions, NeRF family can learn geometrically (i.e., multiview) consistent representations and realize view-dependent rendering results by learning on both position and viewing direction.

**Limitation.** All these methods perform a pixel-wise manner despite architecture difference. Rays (or pixels) are treated independently during training and inference process. The cardinality of the ray set grows quadratically with the increase of image resolution. For an image of 4K ultra high resolution, there exists over 8 million pixels typically presenting richer details and each of which naturally embodies holistic content of a scene in a finer level than the one on a lower resolution image. If directly using such pixel-wise training mechanism for modelling scenes with extremely high-resolution inputs, these methods may struggle with insufficient representational ability for retaining subtle details, even with increased model capacity (shown in the experimental comparison 5.4), which might worsen the issues of lengthy inference with a tremendous MLP or considerable storage cost by using voxel-grid structures with increased volume dimension.

## 3.2. Overall Framework

To extend conventional NeRF methods to achieve high-quality rendering at ultra high resolutions, one straightforward solution is to first train NeRF models for rendering down-sampled outputs and then train parameterized super-resolution on each view to up-sample these outputs to full scale. However, such a solution would result in the artifacts of inconsistent rendering across viewpoints as local patterns captured in the super-resolution stage lack regularization of geometrical consistency (as shown in the ablation study of joint training in 5.5).

We develop a simple yet effective NeRF-4K framework building on the methodology of NeRF-based volume rendering. We first encode geometric information in a lower

3

resolution space through the use of the *View Consistent Encoder* (VC-Encoder for short) module and recover subtle details in a higher resolution (HR) space via the *View Consistent Decoder* (VC-Decoder for short) module. The method aims to boost the representational ability of NeRF-based models on high-frequency details recovery by integrating 3D-aware local features learned in the observations. We further empirically demonstrate that the framework is capable of improving the visual quality of conventional NeRF methods with standard resolution (i.e., 1K), where the encoder and the decoder are learned in the space with the same resolution (as shown in the ablation study of enhancing with local correlation of ray features in 5.5).

### 3.3. View Consistent Encoder

We instantiate the VC-Encoder based on the formulation defined in the DVGO [36], where voxel-grid based representations are learned to encode geometric structure explicitly,

$$(\mathbf{x}, \mathbf{V}) : \left( \mathbb{R}^3, \mathbb{R}^{N_c \times N_x \times N_y \times N_z} \right) \rightarrow \mathbb{R}^{N_c}, \quad (4)$$

where $N_c$ denotes the channel dimension for density ($N_c = 1$) and color modality, respectively. For each sampling point, the density is estimated by trilinear interpolation equipped with a softplus activation, i.e., $\sigma = \delta\left(\text{interp}\left(\mathbf{x}, \mathbf{V}_d\right)\right)$. The colors are estimated with a shallow MLP,

$$\begin{aligned} \mathbf{c} &= f_{\text{MLP}}\left(\text{interp}\left(\mathbf{x}, \mathbf{V}_c\right), \mathbf{x}, \mathbf{d}\right) \\ &= f_{\text{RGB}}\left(g_\theta(\text{interp}(\mathbf{x}, \mathbf{V}_c), \mathbf{x}, \mathbf{d})\right), \end{aligned} \quad (5)$$

where $g_\theta(\cdot)$ extracts volumetric features for color information, and $f_{\text{RGB}}$ denotes the mapping (with one or multiple layers) from the features to RGB images.

We regard $g(\theta; \mathbf{x}, \mathbf{d})$ as the VC-Encoder and the output $\mathbf{g} = g(\theta; \mathbf{x}, \mathbf{d})$ denotes the volumetric feature for the point $\mathbf{x}$ with the viewing direction $\mathbf{d}$, which has embedded geometric information into the feature. In this regard, we can get the descriptor for each ray (or pixel) by accumulating the features of sampling points along the ray $\mathbf{r}$ as in Eqn.1 ,

$$\mathbf{f}(\mathbf{r}) = \sum_{i=1}^{N} T_i \cdot \alpha_i \cdot \mathbf{g}_i. \quad (6)$$

Assume the spatial dimension is $H' \times W'$, the formed feature maps $\mathbf{F}_{\text{en}} \in \mathbb{R}^{C' \times H' \times W'}$ are fed into the VC-Decoder for pursuing high-fidelity reconstruction of fine details.

### 3.4. View Consistent Decoder

For better use of geometric properties embedded in the VC-Encoder, we also generate a depth map $\mathbf{M} \in \mathbb{R}^{H' \times W'}$ by estimating the depth along the camera axis for each ray $\mathbf{r}$,

$$M(\mathbf{r}) = \sum_{i=1}^{N} T_i \cdot \alpha_i \cdot t_i, \quad (7)$$

where $t_i$ denotes the distance of the sampling point $i$ to the camera center as in Eqn.1. The estimated depth map provides a strong guidance for understating the 3D structure of a scene, e.g., nearby pixels on the image plane may be far away in the original 3D space.

VC-Decoder is constructed by taking the feature maps $\mathbf{F}_{\text{en}} \in \mathbb{R}^{C' \times H' \times W'}$ and $\mathbf{M} \in \mathbb{R}^{H' \times W'}$ as input and realizing view synthesis with a higher spatial dimension $H \times W$ through the convolutional neural network $\Psi : (\mathbf{F}_{\text{en}}, \mathbf{M}) \mapsto \mathbf{P}$, where $\mathbf{P} \in \mathbb{R}^{3 \times H \times W}$. $H = sH'$ and $W = sW'$ where $s$ indicates the up-sampling scale. The network is built by stacking several convolutional blocks (with neither non-parametric normalization nor down-sampling operations) interleaved with up-sampling operations. Particularly, instead of simply concatenating the features $\mathbf{F}_{\text{en}}$ and the depth map $\mathbf{M}$, we regard depth signal separately and inject it into every block through a learned transformation to modulate block activations.

Formally, suppose $\mathbf{F}^k$ denotes the activations of an intermediate block with the channel dimension $C_k$. The depth map M passes through the transformation (e.g., with $1 \times 1$ convolution) to predicted scale and bias values with the same channel dimension $C_k$, which are used to modulate $\mathbf{F}^k$ according to:

$$\tilde{\mathbf{F}}_{i,j}^k = \gamma_{i,j}^k(\mathbf{M}) \odot \mathbf{F}_{i,j}^k + \beta_{i,j}(\mathbf{M}). \quad (8)$$

where $\odot$ denotes the element-wise product, $i$ and $j$ indicate the spatial position. More detailed descriptions for the network architecture can be founded in the implementation section and supplemental material.

**Discussion.** Integrating local information of nearby pixels has proven to be effective for recovering high frequency details in single image super-resolution. Pursing high quality super-resolution with view consistency is necessary for the VC-Decoder. Learning local correlation on ray features naturally connects the extraction of spatial patterns to the underlying 3D structure and the modulation induced by depth maps further introduces geometric information to guide the learning.

## 4. Training

The VC-Encoder and the VC-Decoder are jointly trained and the overall framework can be trained in a differentiable and end-to-end manner.

**Patch-based Ray Sampling.** Unlike the pixel-wise mechanism in the traditional NeRF methods, our method aims to capture spatial information between rays (pixels). Therefore, the strategy of random ray sampling is unsuitable here. We present a training strategy with patch-based ray sampling to facilitate the capture of spatial dependencies between ray features.

For training, we first split the images of training views into patches $\mathbf{p}$ with the size $N_p \times N_p$ in order to ensure the

sampling probability on pixels are uniform. When the image spatial dimension can not be exactly divided by the patch size, we truncate the patch until edge. Then we can obtain a set of training patches. A patch (or multiple patches) is randomly sampled from the set, and the rays cast through the pixels in the patch form the min-batch of each iteration.

**Loss Functions.** We found that only using distortion-oriented loss (e.g., MSE, $\ell_1$ and Huber loss) as objective tends to produce blurry or over-smoothed visual effects on fine details. In order to solve the problem, we add the adversarial loss and the perceptual loss to regularize fine detail synthesis. The adversarial loss is calculated on the predicted image patches via the VC-Decoder and training patches through a learnable discriminator which aims distinguish the distribution of training data and predicted one. The perceptual loss $\mathcal{L}_{\mathrm{perc}}$ estimates the similarity between predicted patches $\hat{\mathbf{p}}$ and Ground-Truth $\mathbf{p}$ in the feature space via a pretrained 19-layer VGG network $\varphi$ [34],

$$\mathcal{L}_{\mathrm{perc}} = \|\varphi(\hat{\mathbf{p}}) - \varphi(\mathbf{p})\|_2^2. \tag{9}$$

We use $\ell_1$ loss instead of MSE for supervising the reconstruction of high-frequency details,

$$\mathcal{L}_1 = \frac{1}{N_p^2} \left| \mathbf{C}(\hat{\mathbf{p}}) - \mathbf{C}(\mathbf{p}) \right|. \tag{10}$$

We add an auxiliary MSE loss to facilitate the training of VC-Encoder with the down-scaled images of training views. Formally, the ray features produced by the VC-Encoder are fed into an extra fully-connected layer to regress RGB values in the lower-resolution images. The overall training objective is defined as,

$$\mathcal{L} = \lambda_h \mathcal{L}_1 + \lambda_a \mathcal{L}_{\mathrm{adv}} + \lambda_p \mathcal{L}_{\mathrm{perc}} + \lambda_l \mathcal{L}_{\mathrm{MSE}}^l. \tag{11}$$

where $\lambda_h$, $\lambda_a$, $\lambda_p$ and $\lambda_l$ denote the hyper-parameters for weighting the losses.

The regularization from multiple losses encourage the learning of discriminative patterns in the VC-Decoder, to retain subtle details in 3D scene observations with extremely high resolution. More importantly, joint training can also assist 3D geometry modelling in the encoder through the paths of depth maps and ray features, enabling high-quality view-consistent synthesis across smoothly varying viewpoints.

# 5. Experiments

## 5.1. Implementation

*VC-Encoder Architecture.* We use the configuration of DVGO as the default setting for the encoder. Specially, we extract the ray features at the penultimate layer of the MLP (with the channel dimension 64) with volume rendering in feature spaces following a dimensional reduction layer (with

the channel dimension 6). Then the obtained features are fed into the decoder.

*VC-Decoder Architecture.* We employ a residual skip-connected convolutional blocks [42] for the decoder. Specifically, our decoder consists of a backbone of 5 blocks and an up-sampling head to produce full-scale images. We plug the depth modulation module at the end of each block. For the ablation study of enhancing conventional radiance field methods without resolution change, we exclude the up-sampling head in VC-Decoder. Detailed network architecture can refer to appendix.

*Training.* To facilitate training convergence, in practice we initialize the encoder by pretraining it with 30k epochs following the training setting of DVGO . We then jointly train the encoder and the decoder for 200k epochs. The loss parameters $\lambda_h$, $\lambda_p$, $\lambda_a$ and $\lambda_l$ are respectively set to 1.0, 0.5, 0.02 and 1.0. The learning rates for updating the network weights of the encoder and the decoder are 1e-4 and 2e-4.

## 5.2. Evaluation Metrics

PSNR for evaluating distortion is used as the default metric in NeRF methods, while the metric is insensitive for the artifacts like over-smoothly or blurry details, which has been well-analyzed in [3]. PSNR measures the pixel difference between two images, but it can not measure the perceptual effect of human. In order to evaluate visual quality in a comprehensive manner, we also introduce the LPIPS [51] and NIQE [27] metrics following traditional image enhancement works [5, 42]. Detailed introduction about metrics can refer to appendix.

## 5.3. Datasets

*LLFF.* The LLFF dataset [25] provides the real-world scenes of training views with 4K ultra high resolution. Therefore we use it to conduct experiments and ablation studies by default. It is composed of 8 forward-facing scenes and different scenes have different numbers of training views, between 20 and 60. The original resolution is $4032 \times 3024$ while existing NeRF-based methods use $4\times$ down-scaled images ($1008 \times 756$) for training and inference. In our experiments, we use the original 4K images as groundtruth for training and evaluation in the main experiments. We use corresponding low-dimension ones in the ablation study of assessing the effect of framework for visual quality improvement at different resolutions (i.e., 2K and 1K). We follow other methods to use the camera poses estimated by COLMAP [31].

*Synthetic-NeRF.* The dataset consists of the images rendered from 8 synthetic objects. Each scene is with the fixed 100 training views and the other 200 testing views. As the dataset only provides the resolution of $800 \times 800$, we use it in the ablation study to validate generalization of the method at lower resolution.

Figure 3. **4K visualization effect of comparison methods.** Due to the small number of training data in the fern and leaves scenes, it is difficult for comparison methods to recover high-frequency details, such as the leaf texture of plants. In contrast, our method recovers fine details and outperforms all other methods significantly regarding visual clarity. In T-rex scene, our method can recover clear background railings and even glass shadows, while the baselines show inferior results for rendering the background.

## 5.4. Comparisons

**Quantitative evaluation.** We first conduct the experiments to validate the method on view synthesis at 4K ultra high resolution by comparing it with modern NeRF methods, including Plenoxels [46], DVGO [36], JaxNeRF [9], MipNeRF-360 [2] and NeRF-SR [39]. We also provide inference time and cache memory for reference for a comprehensive evaluation. For a fair comparison with the baselines, we experimented with two settings for them, 1) with standard configuration expect training on 4K resolution, 2) using a large configuration with doubled network parameters and
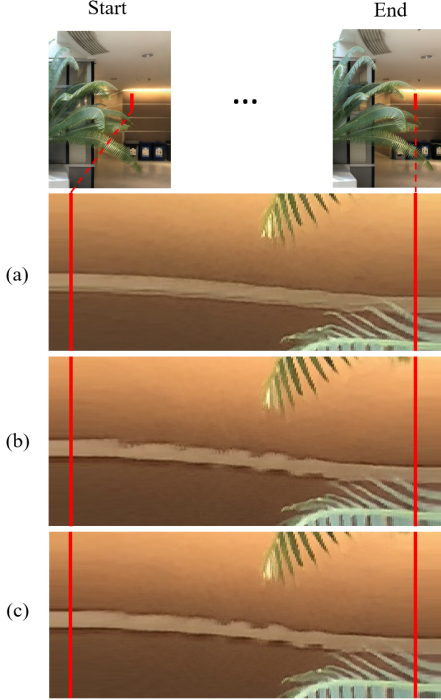
Figure 4. **View consistency visualization**. From horizontal view interpolation videos, we extract a short vertical segment pixel at fixed image location every frame and stack them horizontally to compare view consistency between (a) full training, (b) w/o depth modulation and (c) w/o joint training. Refer to supplemental video for better visualization.

voxel sizes if used and training on 4K resolution. We provide two settings for training our method with GAN-oriented and L1-oriented loss, which means the respective loss contributes more for supervision.

The results are shown in Table 1. Our method achieves obvious advantage in the perception metrics (i.e., LPIPS and NIQE) compared to all the baselines with both settings. The performance on PSNR is also competitive. It is notable that our method can realize higher-fidelity synthesis on preserving fine details compared to DVGO-large (shown in the following qualitative results) although it presents higher PSNR. The observation is also consistent with the results of our method with GAN-oriented and L1-oriented loss, i.e., GAN-oriented training encourages higher performance on visual index as well better qualitative results (shown in Fig. 3), while L1-oriented loss achieve a higher result on PSNR.

Our method achieves compelling performance on both inference efficiency and memory cost, allowing to render an 4K image within 300 ms. Compared to the direct counterpart DVGO, our method realizes significant improvement, i.e., over $20\times$ faster inference with $1/4$ memory overhead.

**Qualitative comparison.** For a better understanding the effect on visual quality, we provide some qualitative comparison in the Fig. 3. The baseline methods show inferior
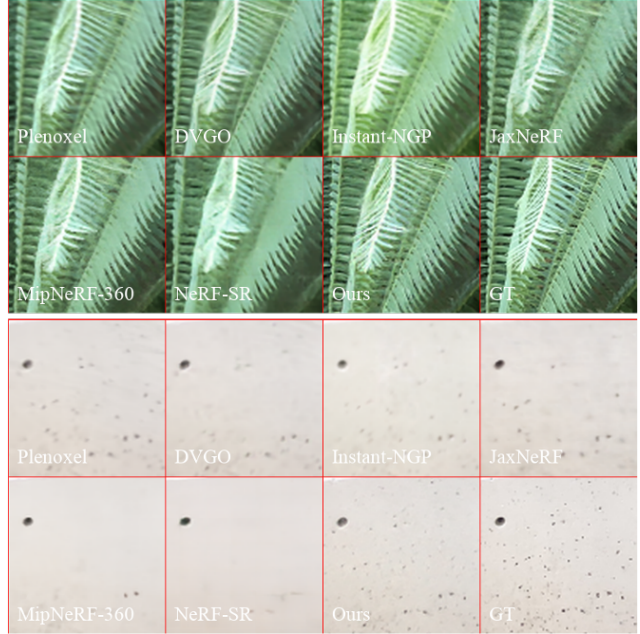


Figure 5. **Visual comparison on default 1K resolution LLFF dataset.** Our methods still achieve the best visual results compared to other methods.

ability on reconstructing subtle details at 4K scenes, incurring details lost or blur. The problem becomes fairly obvious when training views are fewer (e.g., for the scene "fern" and "leaves"). In contrast, our method present high-quality photorealistic rendering on these complex and high-frequency details even on the scenes with limited training views.

## 5.5. Ablation studies

**Enhancing with local correlation of ray features.** The framework can be also adapted to improve the rendering quality of NeRF methods with the standard lower resolution, which is demonstrated in the Table 2 and the Fig. 5. Both quantitative and qualitative results validate the generalization of the framework, i.e., capable of enhancing details recover by integrating local correlation of ray features captured by VC-decoder. It is consistent with As the most commonly used module in image processing, convolution can obtain images' critical local area information with efficient calculation. Previous NeRF was limited to random ray sampling and could not use convolution. Our patch-based ray sampling can combine the rendering process with convolution very well. Fig. 5 shows the comparison of our method with others without changing the resolution. It can be seen that after using convolution, the details of images can be significantly improved even without up-sampling.

**Joint training.** We use patch-based ray sampling in our experiments to jointly train the geometric VC-Encoder and the high-frequency detail VC-Decoder. Experiments

| method | setting | visual metric | | distortion metric | inference time (s) | cache memory (GB) |
| | | LPIPS↓ | NIQE↓ | PSNR↑ | | |
|---|---|---|---|---|---|---|
| Plenoxels | standard | 0.48 | 8.86 | 24.56 | 1.88 | 29.1 |
| | large | 0.48 | 8.85 | 24.57 | 4.02 | 74.0 |
| DVGO | standard | 0.44 | 7.89 | 25.13 | 5.68 | 58.6 |
| | large | 0.39 | 7.05 | **25.53** | 10.39 | 72.6 |
| JaxNeRF | standard | 0.42 | 7.03 | 25.37 | 134.62 | 77.8 |
| | large | 0.39 | 6.80 | 25.50 | 279.83 | 77.8 |
| MipNeRF-360 | standard | 0.37 | 6.31 | 25.34 | 51.38 | 78.1 |
| | large | 0.34 | 5.94 | 25.49 | 105.47 | 78.1 |
| NeRF-SR | standard | 0.52 | 9.26 | 24.15 | 129.19 | 46.7 |
| | large | - | - | - | - | - |
| Ours | GAN | **0.24** | **4.75** | 24.71 | **0.28** | **14.9** |
| | L1 | 0.41 | 7.45 | 25.44 | 0.28 | **14.9** |

Table 1. **Quantitative Comparison on 4K-LLFF dataset.** Quantitative Comparison on 4K-LLFF dataset. We compare with other related works on both visual metric and distortion metric. Our method ranks first on LPIPS and NIQE while having a comparable distortion performance. Moreover, our 4K-NeRF has the fastest inference speed and brings the least memory overhead.
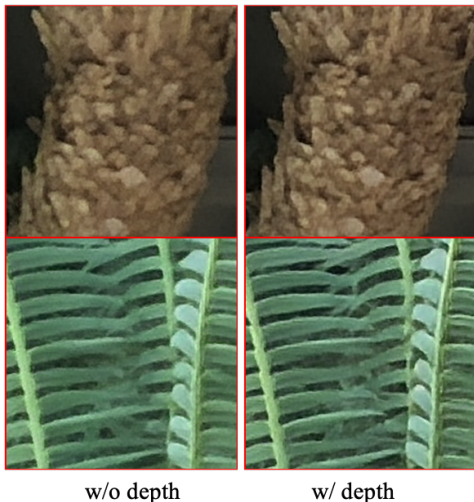


w/o depth        w/ depth

Figure 6. **Ablation of depth modulation.** The clarity of complex and subtle object details can be obviously improved after adding depth modulation.

| Method | LPIPS↓ | NIQE↓ | PSNR↑ |
|---|---|---|---|
| Plenoxels | 0.1345 | 5.0588 | 23.334 |
| DVGO | 0.1672 | 5.7409 | 23.414 |
| JaxNeRF | 0.1850 | 4.3273 | 23.331 |
| MipNeRF-360 | 0.1529 | 3.8579 | 23.466 |
| NeRF-SR | 0.2193 | 4.6946 | 23.232 |
| Ours-L1 | **0.1256** | **3.1597** | **23.869** |

Table 2. Comparison on one of the representative scene (ship) in Synthetic-NeRF dataset.

in Fig. 4. Furthermore, it can improve the details of objects close to the view plane. Fig. 6 shows the change in network recovery details before and after adding depth information.

## 6. Conclusion

In this paper, we explored the NeRF's ability on modelling fine details and proposed a novel framework to boost its representational power on recovering view-consistent subtle details in the scenes with extremely high resolutions. A pair of encoder-decoder modules are introduced to enable modelling geometric properties effectively in a lower resolution space and realizing view consistent enhancement in the full-scale space by virtue of local correlation captured between geometry-aware features. Patch-based sampling training of the framework allow the method integrating the supervision from perception oriented regularization beyond pixel-level comparison. Experiments on challenging real-world datasets validate that our framework is able to realize high fidelity rendering on fine details while keeping view consistency. We expect to investigate the effect of incorporating the framework into the modelling of dynamic scenes as

show that this joint training is essential in maintaining view consistency. We spliced pixel strips of a certain length at a fixed position of each frame in the rendered video. The margin of texture jitter can judge the consistency of objects under different views. Fig. 4 shows texture jitter during spread training, but it is relieved after joint training.

**Depth modulation.** High-resolution details are another critical point for novel view synthesis in ultra-high-resolution scenes, in addition to the importance of view consistency. The near objects are clear, and the distant blur is evident in human vision. We add depth maps as pixel-level modulation information to the VC-Decoder in multi-positions. Experiments show that this approach promotes view consistency

well as neural rendering tasks beyond as a future direction.

# References

[1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *arXiv: Computer Vision and Pattern Recognition*, 2021.

[2] 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*, pages 5470–5479, 2022.

[3] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6228–6237, 2018.

[4] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew DuVall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh representation. *ACM Transactions on Graphics*, 2020.

[5] Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Investigating tradeoffs in real-world video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5962–5971, 2022.

[6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *In Proceedings of the European Conference on Computer Vision*, 2022.

[7] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11065–11074, 2019.

[8] Paul Debevec, Camillo J. Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: a hybrid geometry- and image-based approach. In *ACM SIGGRAPH*, 1996.

[9] Boyang Deng, Jonathan T. Barron, and Pratul P. Srinivasan. JaxNeRF: an efficient JAX implementation of NeRF. https://github.com/google-research/google-research/tree/master/jaxnerf, 2020.

[10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.

[11] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015.

[12] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.

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

[14] Peter Hedman, Suhib Alsisan, Richard Szeliski, and Johannes Kopf. Casual 3d photography. *ACM Transactions on Graphics*, 2017.

[15] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.

[16] Ronghang Hu, Nikhila Ravi, Alexander C Berg, and Deepak Pathak. Worldsheet: Wrapping the world in a 3d sheet for view synthesis from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12528–12537, 2021.

[17] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. Efficientnerf efficient neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12902–12911, 2022.

[18] Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 466–467, 2020.

[19] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1646–1654, 2016.

[20] Andreas Kurz, Thomas Neff, Zhaoyang Lv, Michael Zollhöfer, and Markus Steinberger. Adanerf: Adaptive sampling for real-time rendering of neural radiance fields. In *European Conference on Computer Vision*, pages 254–270. Springer, 2022.

[21] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1833–1844, 2021.

[22] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017.

[23] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *Advances in Neural Information Processing Systems*, 2020.

[24] Nelson L. Max. Optical models for direct volume rendering. *IEEE Trans. Vis. Comput. Graph.*, 1995.

[25] 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*, 2019.

[26] 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 *In Proceedings of the European Conference on Computer Vision*, 2020.

[27] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a "completely blind" image quality analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012.

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

[29] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14335–14345, 2021.

[30] 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*, pages 12892–12901, 2022.

[31] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[32] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8028–8038, 2020.

[33] Horst D Simon. The lanczos algorithm with partial reorthogonalization. *Mathematics of computation*, 42(165):115–142, 1984.

[34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.

[35] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Photo tourism: exploring photo collections in 3d. *ACM SIGGRAPH and Transactions on Graphics*, 2006.

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

[37] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 551–560, 2020.

[38] 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. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2022.

[39] Chen Wang, Xian Wu, Yuan-Chen Guo, Song-Hai Zhang, Yu-Wing Tai, and Shi-Min Hu. Nerf-sr: High quality neural radiance fields using supersampling. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6445–6454, 2022.

[40] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with

pure synthetic data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1905–1914, 2021.

[41] Xintao Wang, Liangbin Xie, Ke Yu, Kelvin C.K. Chan, Chen Change Loy, and Chao Dong. BasicSR: Open source image and video restoration toolbox. https://github.com/XPixelGroup/BasicSR, 2022.

[42] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018.

[43] 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.

[44] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.

[45] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018.

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

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

[48] 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.

[49] Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4791–4800, 2021.

[50] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arxiv CS.CV 2010.07492*, 2020.

[51] 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.

[52] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 286–301, 2018.

# 4K-NeRF: High Fidelity Neural Radiance Fields at Ultra High Resolutions

# Appendix
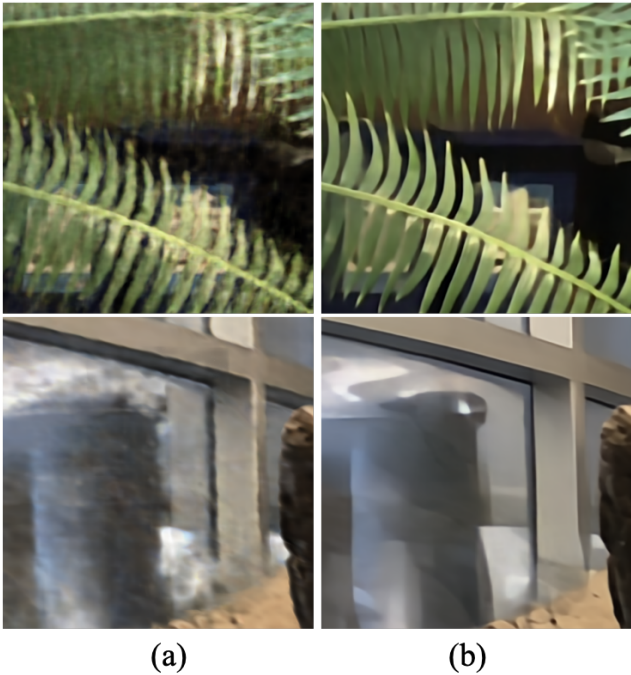


(a)                            (b)

Figure 7. Visual comparison based on TensoRF in 4K scenes. (a) TensoRF. (b) Our results.

## A. Details of Model Structure

We use the default configuration of DVGO [36] as the encoder setting in the experiments. Specifically, the size of voxels is $384 \times 384 \times 256$, and each voxel contains a density value representing geometry and a 12-dimensional color feature followed by a MLP. We extract ray features from the MLP with the channel dimension 64 following a dimensional reduction layer with the channel dimension 6. The encoder is trained on the resolution $1008 \times 756$.

The illustration of the 4K-NeRF structure is shown in Fig. 8. The decoder consists of 5 residual-in-residual dense modules (RRDB) [41, 42] with depth modulation (DM-) as well as one super-resolution head. Each module is comprised of three DM-RRDB blocks interleaved with depth modulation units. We also insert a depth modulation unit for each DM-RRDB block. More detailed configuration can refer to the network configuration provided in the source code. Resolution increase performs in the super-resolution head by stacking two convolutional layers interleaved with $2\times$ bi-linear upsamling operation.

| Scene | Method | LPIPS ↓ | NIQE ↓ | PSNR ↑ |
|-------|--------|---------|--------|--------|
| Fern | TensoRF | 0.464 | 7.172 | 23.333 |
| | Ours | 0.342 | 6.089 | 23.269 |
| Horns | TensoRF | 0.452 | 7.051 | 26.194 |
| | Ours | 0.387 | 6.276 | 26.722 |

Table 3. Quantitative comparison based on TensoRF in 4K scenes.

## B. More Descriptions on Evaluation Metrics

Existing NeRF methods are typically supervised by pixel-level MSE loss and estimated by its direct counterpart PSNR metric. However, only using pixel-level loss is intractable to estimate problems like over-smooth details and blurry visual artifacts. These issues have been well analyzed and explained in detail in the papers [3, 51], revealing the relation between perceptual quality and the degree of distortion. Distortion-oriented metrics (such as PSNR) can be treated as a visual lower bound, ensuring that semantic content in the image is consistent when reaching a certain level. The perceptual effects towards human vision, such as texture details and sharpness, can be measured by virtue of perception-oriented metrics, e.g., LPIPS. PSNR may be inconsistent with visual quality estimated by human eyes. This phenomenon is often more pronounced in ultra-high-resolution videos. Therefore, to quantify and compare the results more reasonably, we use LPIPS and NIQE as evaluation metrics besides PSNR. LPIPS and PSNR are calculated based on test ground-truth views (whose number is limited). As NIQE is a GT-free metric, we calculate across frames of rendered videos given camera trace to better assess cross-view quality.

## C. More Ablation Studies

**Base Architecture.** We use the architecture of DVGO as the base of the VC-Encoder, which can be instantiated with other NeRF-based architectures. In order to further assess the generalization of our framework, we used TensoRF [6] as the VC-Encoder base for training a 4K-NeRF variant and compare the results qualitatively and quantitatively in Fig. 7 and Table. 3. Clear improvements can be achieved on both evaluation metrics and visual qualities. Our 4K-NeRF variant can significantly boost rendering quality on fine details and reduce blurry artifacts even on challenging transparent/translucent objects.

**Loss Functions.** The setting of "L1" denotes training with $\ell_1$ loss only and the setting of "GAN" denotes supervising by adversarial and perceptual losses with $\ell_1$ loss. Besides quantitative results shown in the main paper, we also

Figure 8. The scheme of 4K-NeRF in detail.



(a) Orchids scene in LLFF dataset.



(b) Room scene in LLFF dataset.

Figure 9. Visual comparison of GAN loss and L1 loss on different scenes.

provide visual comparison in the Fig. 9, which shows that using perception-oriented loss helps reducing blurry and over-smooth artifacts and can improve visual quality obviously compared to training with distortion-oriented loss.

## D. Detailed Results

We present detailed results for each 4K scene in the LLFF dataset in the following tables. In addition, we

provide rendered videos on representative scenes ("Fern" and "Horns") for better illustrating the improvement on visual effects achieved by our 4K-NeRF in the https://github.com/frozoul/4K-NeRF, which we recommend to watch on the 4K ultra-high-resolution display.

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|-------|--------|---------|--------|--------|--------|--------|--------|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Fern | Plenoxels | standard | 0.456 | 7.721 | 23.842 | 2.3 | 32.8 |
| | | large | 0.453 | 7.679 | 23.846 | 3.9 | 78.0 |
| | DVGO | standard | 0.424 | 6.910 | 23.741 | 6.2 | 20.1 |
| | | larege | 0.346 | 5.543 | 23.684 | 12.6 | 68.3 |
| | JaxNeRF | standard | 0.399 | 5.623 | 23.470 | 134.7 | 77.8 |
| | | large | 0.354 | 5.312 | 22.689 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.348 | 5.229 | 23.867 | 51.3 | 78.1 |
| | | large | 0.321 | 4.986 | 23.900 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.516 | 7.362 | 22.893 | 129.6 | 46.7 |
| | Ours | GAN | 0.223 | 4.201 | 23.494 | 0.3 | 11.8 |
| | | L1 | 0.353 | 6.377 | 23.691 | 0.3 | 11.8 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|-------|--------|---------|--------|--------|--------|--------|--------|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Flower | Plenoxels | standard | 0.516 | 10.42 | 26.103 | 2.5 | 29.3 |
| | | large | 0.518 | 10.48 | 26.133 | 4.2 | 77.2 |
| | DVGO | standard | 0.500 | 9.964 | 26.857 | 5.6 | 26.5 |
| | | larege | 0.456 | 8.931 | 27.368 | 10.7 | 78.4 |
| | JaxNeRF | standard | 0.489 | 9.308 | 26.783 | 134.7 | 77.8 |
| | | large | 0.458 | 8.890 | 27.265 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.437 | 7.824 | 27.119 | 51.3 | 78.1 |
| | | large | 0.402 | 7.287 | 27.280 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.556 | 11.07 | 25.578 | 129.6 | 46.7 |
| | Ours | GAN | 0.275 | 5.525 | 26.454 | 0.27 | 14.2 |
| | | L1 | 0.493 | 9.514 | 26.865 | 0.27 | 14.2 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|-------|--------|---------|--------|--------|--------|--------|--------|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Fortress | Plenoxels | standard | 0.491 | 9.919 | 28.852 | 2.4 | 30.1 |
| | | large | 0.488 | 9.937 | 28.854 | 4.1 | 76.1 |
| | DVGO | standard | 0.397 | 8.766 | 29.438 | 5.3 | 30.7 |
| | | larege | 0.353 | 8.055 | 30.029 | 9.5 | 71.9 |
| | JaxNeRF | standard | 0.336 | 7.737 | 30.210 | 134.7 | 77.8 |
| | | large | 0.334 | 7.490 | 29.882 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.314 | 7.472 | 30.169 | 51.3 | 78.1 |
| | | large | 0.294 | 6.806 | 30.231 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.517 | 9.637 | 28.719 | 129.6 | 46.7 |
| | Ours | GAN | 0.227 | 4.857 | 28.120 | 0.25 | 15.3 |
| | | L1 | 0.404 | 8.320 | 29.853 | 0.25 | 15.3 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|---|---|---|---|---|---|---|---|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Horns | Plenoxels | standard | 0.510 | 8.298 | 24.743 | 2.3 | 31.0 |
| | | large | 0.511 | 8.235 | 24.763 | 4.4 | 78.4 |
| | DVGO | standard | 0.462 | 7.053 | 25.632 | 5.4 | 40.8 |
| | | larege | 0.394 | 6.340 | 26.367 | 9.6 | 72.0 |
| | JaxNeRF | standard | 0.430 | 5.945 | 26.127 | 134.7 | 77.8 |
| | | large | 0.402 | 5.853 | 26.760 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.371 | 5.172 | 26.220 | 51.3 | 78.1 |
| | | large | 0.344 | 4.952 | 26.224 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.553 | 9.758 | 23.694 | 129.6 | 46.7 |
| | Ours | GAN | 0.261 | 4.439 | 25.066 | 0.29 | 18.8 |
| | | L1 | 0.399 | 6.241 | 26.336 | 0.29 | 18.8 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|---|---|---|---|---|---|---|---|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Leaves | Plenoxels | standard | 0.520 | 7.749 | 20.028 | 1.0 | 23.3 |
| | | large | 0.518 | 7.712 | 20.030 | 1.5 | 59.1 |
| | DVGO | standard | 0.511 | 7.388 | 20.220 | 5.7 | 22.6 |
| | | larege | 0.447 | 6.568 | 20.211 | 10.1 | 71.0 |
| | JaxNeRF | standard | 0.536 | 6.942 | 19.781 | 134.7 | 77.8 |
| | | large | 0.453 | 6.294 | 20.148 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.427 | 6.078 | 19.835 | 51.3 | 78.1 |
| | | large | 0.379 | 5.593 | 19.970 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.559 | 8.167 | 19.033 | 129.6 | 46.7 |
| | Ours | GAN | 0.267 | 4.367 | 19.781 | 0.25 | 13.4 |
| | | L1 | 0.461 | 7.075 | 19.819 | 0.25 | 13.4 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|---|---|---|---|---|---|---|---|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Orchids | Plenoxels | standard | 0.575 | 9.150 | 19.874 | 2.1 | 35.3 |
| | | large | 0.572 | 9.036 | 19.870 | 3.7 | 69.4 |
| | DVGO | standard | 0.539 | 8.112 | 20.098 | 6.1 | 22.5 |
| | | larege | 0.491 | 6.924 | 19.970 | 14.3 | 73.7 |
| | JaxNeRF | standard | 0.549 | 7.872 | 19.649 | 134.7 | 77.8 |
| | | large | 0.498 | 7.147 | 19.383 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.482 | 6.880 | 19.511 | 51.3 | 78.1 |
| | | large | 0.432 | 6.234 | 19.672 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.594 | 8.973 | 19.432 | 129.6 | 46.7 |
| | Ours | GAN | 0.307 | 5.203 | 20.005 | 0.32 | 12.5 |
| | | L1 | 0.523 | 7.649 | 19.557 | 0.32 | 12.5 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|---|---|---|---|---|---|---|---|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| Room | Plenoxels | standard | 0.392 | 9.038 | 28.133 | 2.7 | 20.5 |
| | | large | 0.390 | 9.038 | 28.133 | 5.2 | 75.3 |
| | DVGO | standard | 0.357 | 7.979 | 29.554 | 5.4 | 30.3 |
| | | larege | 0.323 | 7.486 | 30.834 | 9.4 | 71.4 |
| | JaxNeRF | standard | 0.317 | 6.744 | 31.114 | 134.7 | 77.8 |
| | | large | 0.309 | 7.360 | 31.733 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.309 | 6.614 | 30.687 | 51.3 | 78.1 |
| | | large | 0.289 | 6.645 | 31.540 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.407 | 9.316 | 29.369 | 129.6 | 46.7 |
| | Ours | GAN | 0.198 | 4.724 | 29.620 | 0.27 | 15.3 |
| | | L1 | 0.304 | 7.732 | 31.147 | 0.27 | 15.3 |

| Scene | Method | Setting | Visual index | | Distortion index | Inference time | Cache memory |
|---|---|---|---|---|---|---|---|
| | | | LPIPS ↓ | NIQE ↓ | PSNR ↑ | (s) ↓ | (GB) ↓ |
| T-rex | Plenoxels | standard | 0.409 | 8.584 | 24.896 | 2.8 | 30.8 |
| | | large | 0.410 | 8.647 | 24.926 | 5.4 | 78.7 |
| | DVGO | standard | 0.356 | 6.985 | 25.512 | 5.7 | 37.4 |
| | | larege | 0.315 | 6.543 | 25.764 | 10.2 | 79.1 |
| | JaxNeRF | standard | 0.335 | 6.030 | 25.839 | 134.7 | 77.8 |
| | | large | 0.316 | 6.092 | 26.147 | 279.3 | 77.8 |
| | MipNeRF-360 | standard | 0.296 | 5.176 | 25.312 | 51.3 | 78.1 |
| | | large | 0.274 | 5.033 | 25.122 | 105.2 | 78.1 |
| | NeRF-SR | standard | 0.454 | 9.857 | 24.230 | 129.6 | 46.7 |
| | Ours | GAN | 0.185 | 4.672 | 25.121 | 0.26 | 18.0 |
| | | L1 | 0.324 | 6.716 | 26.276 | 0.26 | 18.0 |