

Neural Volume Super-Resolution

Yuval Bahat¹Yuxuan Zhang¹Hendrik Sommerhoff²Andreas Kolb²Felix Heide¹¹Princeton University ²University of Siegen

Abstract

Neural volumetric representations have become a widely adopted model for radiance fields in 3D scenes. These representations are fully implicit or hybrid function approximators of the instantaneous volumetric radiance in a scene, which are typically learned from multi-view captures of the scene. We investigate the new task of neural volume super-resolution – rendering high-resolution views corresponding to a scene captured at low resolution. To this end, we propose a neural super-resolution network that operates directly on the volumetric representation of the scene. This approach allows us to exploit an advantage of operating in the volumetric domain, namely the ability to guarantee consistent super-resolution across different viewing directions. To realize our method, we devise a novel 3D representation that hinges on multiple 2D feature planes. This allows us to super-resolve the 3D scene representation by applying 2D convolutional networks on the 2D feature planes. We validate the proposed method by super-resolving multi-view consistent views on a diverse set of unseen 3D scenes, confirming qualitative and quantitatively favorable quality over existing approaches.

1. Introduction

Reconstructing latent high-quality images from images captured under non-ideal imaging conditions is a broadly studied research field. A large body of existing work explores methods for removing image noise [8, 21], sharpening blurry images [40, 2], and increasing image resolution [30, 22]. This direction includes recent works that are based on deep learning [22, 40], mapping degraded images to their latent non-degraded counterparts, and have led to a significant leap in performance. At the heart of these methods typically lies the application of a series of 2D convolution operations to the degraded image [25, 22], which eventually yields the reconstructed image at a higher resolution.

Concurrently with research on image reconstruction, a rapidly growing body of work explores neural rendering, allowing not only to reconstruct a single image but unseen novel views from a set of observed images. In their work

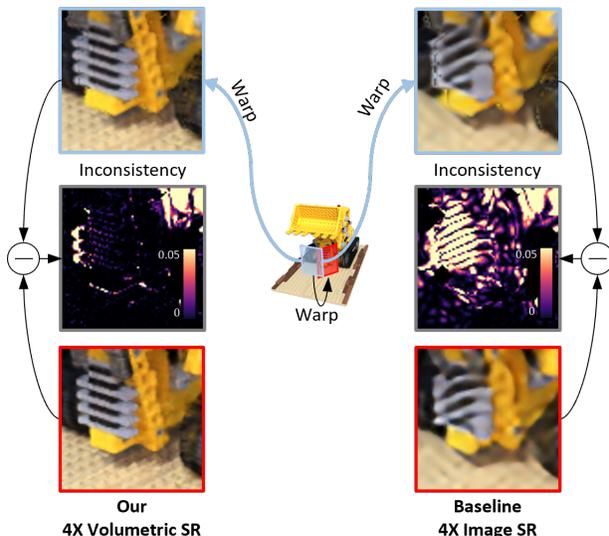


Figure 1. **Geometrically consistent super-resolution.** Our method (bottom left) renders HR views of a scene captured in LR by super-resolving its volumetric representation. This guarantees the geometric consistency of the rendered scene across viewing directions, and prevents inconsistent distortions such as the slanted bulldozer grill in the output of the EDSR [25] image SR method (bottom right). Here we visualize the across-views inconsistency by rendering two adjacent high-resolution scene views, and then warping one of them (top) onto the perspective of the other (bottom). We then compute error maps between each rendered pair (middle row), which further highlights our enhanced consistency.

on Neural Radiance Fields (NeRF) [32] in 2020, Mildenhall *et al.* showed that state-of-the-art rendering quality can be achieved even with a lean Multilayer Perceptron (MLP) architecture as a coordinate-based representation for a 5D radiance field. Many works have since proposed improved and more efficient ways to represent 3D scenes based on a given finite set of 2D images, including methods that tackle modified lighting conditions [36] or support dynamic scenes [33]. These existing methods directly supervise the scene representation using the observed images.

In this work, we explore an alternative direction at the intersection of image reconstruction and neural rendering,

and propose the new task of neural volume super-resolution - learning to super-resolve neural volumes. Specifically, given a finite set of low-resolution images corresponding to a 3D scene, we aim to render high-resolution images of the scene corresponding to arbitrary possible viewing directions, including novel unseen ones. Rather than operating on the rendered (2D) images and upsampling these images after generation, we propose a novel approach to tackle this task, by operating directly on the volumetric (3D) scene representation. An advantage of this approach, illustrated in Fig. 1, is that it guarantees by design a *geometrically consistent* image reconstruction across all viewing directions, in contrast to operating in the 2D domain and separately reconstructing each rendered image.

As a practical and efficient approach for enabling upsampling neural volumes, we propose a new 2D plane-based scene representation that allows us to employ 2D convolutional neural networks. Specifically, we represent a 3D scene by implicitly associating each point in the 3D volume with a density value and a (view-direction dependent) RGB value following existing radiance fields works, while investigating an alternative way to infer these values. Chan *et al.* [9] accumulated features for each 3D point by projecting its location onto three 2D feature planes, for representing human faces. We instead propose a plane representation that allows representing general scenes, and introduce a fourth feature plane to enable view direction dependency. Given a desired viewing direction, we can reconstruct each image pixel by marching along its corresponding ray from the camera, taking into account the cumulative effect of the density and RGB values predicted by our model for each point along the ray.

Rendering a *super-resolved* image from the proposed scene model, that was learned from a set of *low-resolution* images, requires a single modification to our pipeline, while the rest remains unchanged: Before projecting each point in the volume onto the 2D planes to infer its density and RGB values, we super-resolve the three 2D planes using a learned super-resolution module that operates in feature space. This approach allows us to take advantage of existing methods explored for 2D super-resolution over the years.

We train our *Neural Volume Super-Resolution* (NVSR) framework end-to-end using low-resolution and matching high-resolution image pairs corresponding to synthetic scenes from freely available online datasets. We then apply the trained SR model to the 3D representations of unseen scenes for which only low-resolution images are available. We validate that unlike super-resolving images post-rendering, the resulting super-resolved rendered images are indeed consistent across different views. The proposed method outperforms such conventional post-rendering SR methods both qualitatively and quantitatively, exceeding state-of-the-art image SR methods by over 1 dB in PSNR

when trained on the same dataset.

Specifically, we make the following contributions:

- We propose a method for super-resolving volumetric neural scene representations, which allows generating geometrically consistent, high-resolution novel views corresponding to a scene captured in low resolution.
- We introduce a super-resolution approach that operates directly in feature space. To this end, we develop a novel feature-plane representation model that allows for rendering viewing-direction dependent views of general 3D scenes.
- We assess our method for super-resolving unseen scenes both qualitatively and quantitatively, validating that it produces sharp, multi-view consistent super-resolved images.

We will release all code and models needed to reproduce the results presented in this work.

2. Related Work

This work introduces a method for reconstructing images of a 3D scene by operating directly on its volumetric neural representation. Next, we provide background on the two most related fields, image super-resolution and neural 3D scene representation.

Image Super Resolution Image super resolution (SR) belongs to a wide category of ill-posed image reconstruction tasks such as image denoising [8, 21], dehazing [23, 3], deblurring [40, 2], inpainting [7, 19] and more. Image SR aims to compensate for information lost during capturing as a result of the lossy image acquisition process, due to a resolution loss in the optical system and finite sampling by the sensor. Existing image SR methods produce a reconstructed image of the captured scene as if captured in ideal settings. To recover the fine details lost due to the finite capturing resolution, some SR works propose to collect additional information from multiple captures of the same scene [12, 42], which typically allows increasing the resolution by only up to a factor of two. To allow higher SR factors, classical methods have proposed to employ various image priors that exploit the unique characteristics of natural images, *e.g.* their internal self-similarity [15, 30], heavy-tailed gradient magnitude distribution [37] and distinctive image edge statistics [38].

Learning Image Super Resolution With the arrival of deep learning, researchers proposed to learn these image characteristics and train Convolutional Neural Networks (CNNs) on large image datasets [10, 25, 24], leading to unprecedented performance in terms of minimizing reconstruction error (*i.e.* increasing PSNR). Following the introduction of Generative Adversarial Networks (GANs) [16],

some methods also began targeting alternative objectives, such as perceptual reconstruction quality [22, 41] or exploring the diverse space of SR solutions [4, 28]. A further line of work on image super resolution takes relevant reference images as additional user input [44, 11] to better adapt to each specific scene. However, common to all of these methods is that they aim to super-resolve a *single* image at a time. To the best of our knowledge, the proposed is the first work aiming to *consistently* super-resolve images corresponding to the *infinite* possible views of a given 3D scene, by operating directly on its neural representation. While methods for spatial Video SR [20, 26, 14] also attempt to simultaneously handle sequences of inter-related frames, these methods are limited to positions along the original camera trajectory and cannot render novel views. Moreover, they cannot inherently guarantee consistent reconstruction across frames, unlike our approach.

Neural Scene Representations. An increasingly large body of work learns representations of 3D scenes, aiming to allow synthesizing novel scene views from arbitrary directions, as well as editing the scene *e.g.* to allow modifying the perceived illumination settings or object properties. Early works [1, 35] tried to learn an explicit 3D scene model by fitting it to a given set of images capturing the scene from different view directions. The recent introduction of Neural Radiance Fields [32] (NeRF) then brought upon a significant leap in reconstructed image quality, rendering photo-realistic views which include fine, high-frequency details. They rely on a lean coordinate-based Multilayer Perceptron (MLP) to implicitly represent the 5D radiance field, made possible mainly thanks to artificially introducing high frequencies into the model input, in the form of sinusoidal positional encoding [39].

In their follow up Mip-NeRF [5, 6] papers, Barron *et al.* focused on handling diverse scene capture resolutions. By modeling sampled rays as conical frustums they were able to reduce blur and aliasing in rendered outputs. In contrast to their work, ours allows rendering novel high-resolution views corresponding to a scene which was originally captured in low-resolution, by *reconstructing* the high-frequency content that was lost during scene acquisition.

Other methods proposed to use alternative implicit representation models based on, *e.g.*, Voxels [27] or spherical harmonics [13]. Chan *et al.* [9] proposed to generate geometrically consistent human faces by representing the volume using three 2D feature planes coupled with a simple MLP decoder. We also employ 2D feature planes to support general 3D scenes and view-direction dependency and take advantage of the planar representation to perform our volumetric SR in 2D, while leveraging the existing rich knowledge base which already exists for 2D SR (see Sec. 3.1).

Many recent methods focused on enabling NeRF to learn

from imperfect image sets, *e.g.* using partly occluded [29] or very few scene captures [18]. In contrast, our work aims to endow the learned representation model with the ability to reconstruct finer details than the ones captured in the given image sets, by learning volumetric scene priors, *i.e.*, a volumetric extension of (natural) image priors.

Neural Radiance Fields for Image Reconstruction. In a separate line of work, some methods harness the implicit 3D representation of NeRF as a tool for image denoising [34] or high dynamic range (HDR) reconstruction [31, 17]. These methods take as input a burst of noisy or low dynamic range scene captures, respectively, and distill the multi-frame signal while relying on the geometrical consistency inherent to the 3D representation. In contrast, our approach allows using a set of degraded (low-resolution) scene captures for synthesizing consistent novel scene views, containing fine details that are *unavailable* in the input images, by relying on *externally learned* volumetric priors (see Sec. 3.2).

3. Super-Resolving Neural Volumes

We introduce a method for super-resolving the representation of a 3D scene s given \mathcal{I}_{lr}^s , a set of low-resolution (LR) captures of the scene and their corresponding relative camera poses. To this end, we devise two principle modules, namely a scene representation model G and a super-resolution model F . We first use the images in \mathcal{I}_{lr}^s to learn the parameters θ_G^s of representation model G , which comprises a decoder model D and learned feature planes \mathcal{P}_s corresponding to each scene s . Once trained, one can use G to render novel low-resolution views of the scene for new (unseen) camera poses:

$$\hat{i}_{lr}(r) = G(\theta_G^s, r). \quad (1)$$

We use r here to denote a pixel in the rendered image and its corresponding camera ray interchangeably, and denote by \hat{i}_{lr} the rendered images, as they too lack the high-resolution fine details which are missing from the images in \mathcal{I}_{lr}^s . To render high-resolution (HR) views containing fine details and textures, we propose to apply model F on the *parameters* of model G which constitute the volumetric representation of s , that is

$$\hat{i}_{sr}(r) = G(F(\theta_F, \theta_G^s), r). \quad (2)$$

Here \hat{i}_{sr} is the rendered super-resolved image and θ_F are the parameters of model F , which are learned independently of any one specific scene.

An overview of our approach is illustrated in Fig. 2. We next elaborate on each of its two modules; we first describe our novel plane-based neural 3D representation model G , and then explain how we use this representation to perform volumetric SR using F .

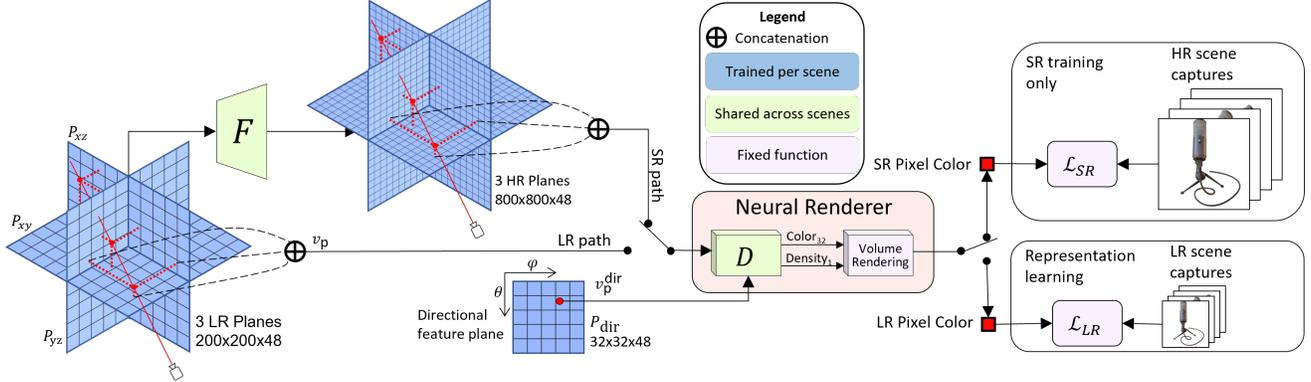


Figure 2. **Overview of our Volumetric SR framework.** We use three LR positional feature planes (left) and one view-direction feature plane (bottom) to represent each 3D scene based on its captured LR images. Image pixels are rendered by projecting points along the corresponding camera ray onto the four planes and extracting the corresponding feature vectors v_p . These are then processed by a decoder MLP D (shared across all scenes) to yield volumetric density and radiance values, which translate into pixel values through volume rendering. To render HR scene views, we first super-resolve the positional feature planes using F , and then extract v_p from the super-resolved planes (top). We use pairs of GT LR-HR image captures from a set of training scenes to train our framework end to end.

3.1. Quadri-Plane Radiance Fields for 3D Scenes

To take advantage of advances in (2D) image SR, and to allow computationally efficient volumetric SR, we propose to represent a scene s with a novel *quadri-plane* model comprised of four multi-channel 2D feature planes of size $N \times N \times C$, coupled with a small decoder neural network.

Similar to existing neural rendering methods [32], we view the volume as a pair of 3D fields representing scene density σ and corresponding emitted radiance (in RGB). We then diverge from the typical scheme and propose a feature plane-based representation. Chan *et al.* [9] recently proposed a scene representation using feature planes for generating human faces. Specifically, they employ explicit feature vectors arranged in three orthogonal axis-aligned feature planes $\{P_{xy}, P_{xz}, P_{yz}\}$. Then each 3D position $p \in \mathbb{R}^3$ is queried by projecting it onto each of the three feature

planes and retrieving the per-plane feature vector corresponding to the point via bilinear interpolation, then aggregating them into a point feature representation vector v_p . This vector is fed into a light-weight *decoder* MLP D to yield the density and RGB values corresponding to p , which are then used to render an image of the scene using neural volume rendering [32]. In this work, we propose several modifications to this approach which are critical for the proposed method to function. We describe these modifications in the following.

View-direction Dependent Radiance. We introduce a 4th feature plane P_{dir} to accommodate effects like non-uniform scene illumination, as we show in Fig. 3. The two axes of this plane reflect the 2D space of possible viewing directions for each point p . We then infer a complementary feature vector v_p^{dir} using bilinear interpolation over P_{dir} at coordinates (θ_p, ϕ_p) , corresponding to the viewing azimuth and elevation of point p , respectively. Features from this plane only affect the radiance and not the density output of D , as unlike radiance, volume density does not vary with viewing direction in the real world [32].

Representing General Scenes. To allow representing general scenes with arbitrary geometric structures beyond the human faces demonstrated in [9], we follow the practice from [32] and use a pair of coarse and fine decoders D_c and D_f , respectively, that are *shared* across all scenes; we first apply the coarse decoder to a set of stratified sample points p along rays corresponding to each pixel. We then use the resulting values to bias a subsequent sampling of points along each ray towards more relevant parts of the volume, and process those using the fine decoder.

Given a set of training scenes \mathcal{S}_t and corresponding sets of captured images, we learn a set of four planes per

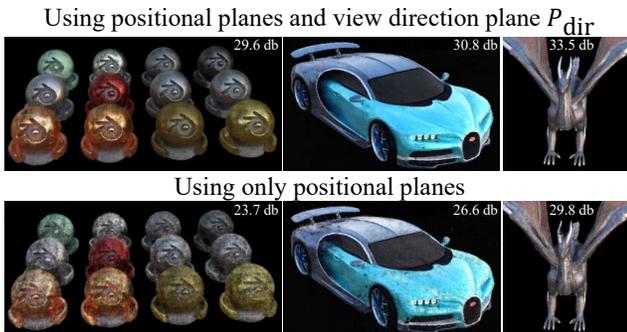


Figure 3. **View-direction rendering dependency.** Utilizing a 4th feature plane P_{dir} (top row) corresponding to viewing direction significantly improves reconstruction performance, as it accounts for directional appearance changes stemming from non-uniform lighting conditions, especially in cases of shiny, reflective objects.

scene $\mathcal{P}_s = \{P_{xy}, P_{xz}, P_{yz}, P_{\text{dir}}\}$, $s \in \mathcal{S}_t$ jointly with the decoder pair parameters θ_D (shared across all scenes), which together constitute the representation model parameters $\theta_G^s = \theta_D \cup \mathcal{P}_s$. To this end, we minimize the following image rendering penalty:

$$\mathcal{L}_{\text{LR}} = \sum_{s \in \mathcal{S}_t, r \in \mathcal{I}_{\text{lr}}^s} \|G_c(\theta_G^s, r) - I_{\text{lr}}(r)\|_2^2 + \|G_f(\theta_G^s, r) - I_{\text{lr}}(r)\|_2^2. \quad (3)$$

Here $I_{\text{lr}}(r)$ is the ground-truth RGB pixel value corresponding to ray r sampled from the training image $I_{\text{lr}} \in \mathcal{I}_{\text{lr}}^s$, and $G_c(r)$ and $G_f(r)$ are the corresponding rendered values obtained through volume rendering over the outputs of the coarse and fine decoders D_c and D_f , respectively.

We tailor the resolution N of the learned planes to the resolution of the images in $\mathcal{I}_{\text{lr}}^s$: As we show in Fig. 4, using planes that are too small prohibits the representation of high visual frequencies, similar to using NeRF [32] without positional encoding. In contrast, using too large planes results in some plane regions having too little relevant information in the captured scene images $\mathcal{I}_{\text{lr}}^s$, which is manifested as visual artifacts when rendering novel scene views. We will revisit this trade-off in plane resolution later in the context of super-resolving the scene.

Note also that we choose a smaller resolution value $N_{\text{dir}} < N$ for the view-direction feature plane P_{dir} , since the image set contains less information relevant to this plane: Each image pixel carries information relevant to all points along its corresponding ray r which, in turn, makes it relevant to many different locations on the positional feature planes $\{P_{xy}, P_{xz}, P_{yz}\}$. In contrast, all points along a ray correspond to the same viewing direction, meaning that the information in each pixel is relevant to a single point on P_{dir} .

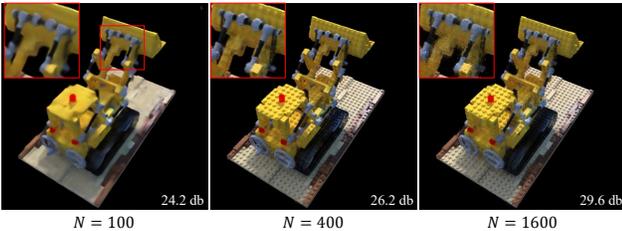


Figure 4. **Effect of feature planes resolution.** Novel view reconstructions using different positional plane resolutions N , all learned from the same set of 400×400 pixel training images. Feature planes that are too small cannot represent fine details (left), while larger planes yield sharp, detailed rendered images (middle). However, at some point increasing the feature plane size leads to visual artifacts (right), as some regions in the oversized feature plane cannot be learned due to the limited resolution of the training images.

3.2. Super-resolving Implicit Representations

To produce HR views of a scene given only a set of its LR captures, one could take a classic super-resolution approach and apply any one of many existing image SR methods on the rendered LR scene views. However, since such methods have no notion of the underlying three-dimensional scene, novel views contain inconsistent super-resolved details in the synthesized views (see Fig. 1). Our approach addresses this issue and allows producing super-resolved, *geometrically consistent* scene views.

As we discuss in Sec. 3.1, employing larger plane resolution N allows us to represent and render higher visual spatial frequencies, corresponding to small details and fine textures. Therefore, substituting the set of planes \mathcal{P}_s learned from the LR scene captures $\mathcal{I}_{\text{lr}}^s$ with a set of higher-resolution planes corresponding to the same scene s can allow for rendering the desired HR views.

We propose to approximate sets of HR planes by learning plane super-resolution model F , which we apply on the learned LR planes (Eq. 2) to obtain α times super-resolved versions of them, with α being the desired SR factor. We use a single model to super-resolve all positional¹ feature planes $P \in \mathcal{P}_s \setminus \{P_{\text{dir}}\}$, while taking advantage of our 2D planes-based representation model by employing an existing image SR network architecture for F . Specifically, we borrow the residual network architecture used for the successful 2D super-resolution method EDSR [25]. To learn the mapping from sets of LR planes to their HR counterparts, we train F on a training set of scenes for which ground truth (GT) HR captures are available, by substituting $G_f(r)$ in Eq. 3 with the output from Eq. 2 and using HR scene captures $I_{\text{hr}} \in \mathcal{I}_{\text{hr}}^s$, resulting in a second loss term:

$$\mathcal{L}_{\text{HR}} = \sum_{s \in \mathcal{S}_t, r \in \mathcal{I}_{\text{hr}}^s} \|G_f(F(\theta_F, \theta_G^s), r) - I_{\text{hr}}(r)\|_2^2. \quad (4)$$

Generalizing to Unseen Scenes. When we trained our SR model F through minimizing \mathcal{L}_{HR} over training scenes in \mathcal{S}_t we observed a severe overfitting problem, which we attribute to the following subtle point: While the input to decoder D corresponding to 3D point p comprises *all three* feature vectors v_p , model F acts by super-resolving each feature plane *independently*. It is trained by minimizing the loss over training scenes in \mathcal{S}_t and their corresponding v_p triplet combinations. This, however, results in poor performance when tested on unseen scenes $s \in \mathcal{S}_e = \{s | s \notin \mathcal{S}_t\}$, consisting of different combinations of v_p which were often not seen during training.

To circumvent this problem we additionally train F on v_p triplets corresponding to test scenes. Since GT HR captures are unavailable for these scenes, we instead exploit

¹We use the viewing direction plane as is without super-resolving it.

their available LR captures, by introducing a downsampling inconsistency loss term $\mathcal{L}_{\text{incon}}$. This term penalizes for the difference between the available GT LR views and their corresponding downsampled HR views rendered by our model:

$$\mathcal{L}_{\text{incon}} = \sum_{s \in \mathcal{S}_e, r \in \mathcal{I}_{\text{lr}}^s} \|[G_f(F(\theta_F, \theta_G^s), r)] \downarrow_{\alpha} - I_{\text{hr}}(r)\|_2^2, \quad (5)$$

Where $[\cdot] \downarrow_{\alpha}$ corresponds to image downsampling by a factor of α and \mathcal{S}_e denotes the set of test scenes.

We train our framework *end-to-end* and learn the parameters of G (including fine and coarse decoder parameters θ_D and LR scene planes \mathcal{P}_s corresponding to each training scene s) together with training our SR model F . To this end, we alternate between the following three optimization steps: (i) Minimizing \mathcal{L}_{LR} (3) over LR training scene images $\{\mathcal{I}_{\text{lr}}^s\}_{s \in \mathcal{S}_t}$, (ii) minimizing \mathcal{L}_{HR} (4) over HR training scene images $\{\mathcal{I}_{\text{hr}}^s\}_{s \in \mathcal{S}_t}$, or (iii) minimizing $\mathcal{L}_{\text{incon}}$ (5) over LR test scene images $\{\mathcal{I}_{\text{lr}}^s\}_{s \in \mathcal{S}_e}$. At each training step we randomly choose one of the three at probability ratios of $f_{\text{LR}} : f_{\text{HR}} : f_{\text{incon}}$, respectively.

3.3. Implementation Details

In both of our decoder models D_c and D_f , we calculate the density value σ corresponding to each 3D point p using fully-connected, 128 channel, 4 layer MLPs, which take as input the average of the positional feature vectors v_p . For computing the corresponding radiance values $\{r, g, b\}$ we feed the concatenation of all 4 (including the view-direction) feature vectors v_p into MLPs with identical characteristics. We use $C = 48$ channels and a spatial resolution of $N = 200$ for our low-resolution positional feature planes learned from scenes captured at 100 pixels, and set the view-direction feature plane P_{dir} resolution to $N_{\text{dir}} = 32$. For our SR model F we use the architecture from EDSR [25] with 32 residual blocks and 256 channels.

We train our entire model end-to-end using Adam optimizers with learning rates of 0.0005 for the decoder and feature planes and 0.00005 for F . At each step, we sample 4096 rays r from an image I corresponding to a scene s , all chosen randomly. We use an NVIDIA A100 GPU, and first train our framework to convergence by minimizing \mathcal{L}_{LR} and \mathcal{L}_{HR} over training scenes $s \in \mathcal{S}_t$ ($f_{\text{LR}} = f_{\text{HR}} = 1$, $f_{\text{incon}} = 0$). Given a test scene $s \in \mathcal{S}_e$, we resume training by alternating between the three optimization steps at a ratio of $f_{\text{LR}} = f_{\text{HR}} = 1$, $f_{\text{incon}} = 10$.

4. Experiments

We next demonstrate the advantage of using our volumetric approach for super-resolving synthetic, as well as real 3D scenes captured in LR. To this end, we use the sets of eight synthetic scenes and eight real scenes available on the NeRF [32] project page, as well as additional 21 scenes

rendered from 3D Blender models available online². These scenes depict complex objects with fine details and textures, which are crucial for demonstrating and evaluating the SR quality. For all synthetic scenes, we use the 100 available training images with dimensions 400×400 as the ground truth HR captures $\mathcal{I}_{\text{hr}}^s$, and then downsample each image by a factor of 4 using the default bicubic kernel to simulate the set of LR captures corresponding to the scene, $\mathcal{I}_{\text{lr}}^s$. Each of our real world scenes has a couple of dozen available images with varying dimensions, which we resize to obtain similar dimensions as in the synthetic case.

For quantitatively evaluating on synthetic scenes, we train our model (including both representation model G and SR model F) on 25 scenes, while holding out the remaining four (‘mic’, ‘ship’, ‘chair’ and ‘lego’) scenes for evaluation. We evaluate all methods on 200 test images corresponding to novel views of each of these unseen scenes, while comparing our approach with methods and baselines representing several different approaches, by using a set of evaluation metrics. We use the same evaluation protocol with real scenes, using six scenes for training and the remaining unseen two (‘fern’ and ‘orchids’) for evaluation. We report the results on synthetic and real scenes separately in Tabs. 1 and 2, respectively, and elaborate next on the metrics we employ, as well as on the baselines we compare to.

Evaluation Metrics. We use PSNR and SSIM to measure how well methods minimize the error with respect to the corresponding GT HR images. To evaluate perceptual quality, we employ the LPIPS [43] score, where lower means better.

Baselines. We compare our approach to several baseline alternatives. We first compare with the naive approach of using models trained on the LR image sets (which lack HR content) to render HR novel views, without attempting to reconstruct the missing HR content. We repeat this experiment twice, once using our quadri-plane model (‘Planes’) and a second time using the Mip-NeRF [5] method (‘Mip-NeRF’), which was designed to handle multiple image resolutions for reducing aliasing and blur.

We additionally compare with alternative ways to reconstruct the missing HR content, by harnessing existing methods for 2D super resolution as pre-processing or post processing steps: In the first (denoted ‘pre. 2D SR’) we apply a 2D SR method on the LR training images and use the resulting super-resolved images to learn a Mip-NeRF model for the scene, while in the second (denoted ‘post. 2D SR’) we use the 2D SR method to super-resolve LR images rendered by a Mip-NeRF model which was learned from the LR images. For the 2D SR step we use the state-of-the-art SwinIR method [24] which employs visual transformers. In

²Links to the these freely available Blender models are provided in the Supp. material. Our entire project code will become available online.

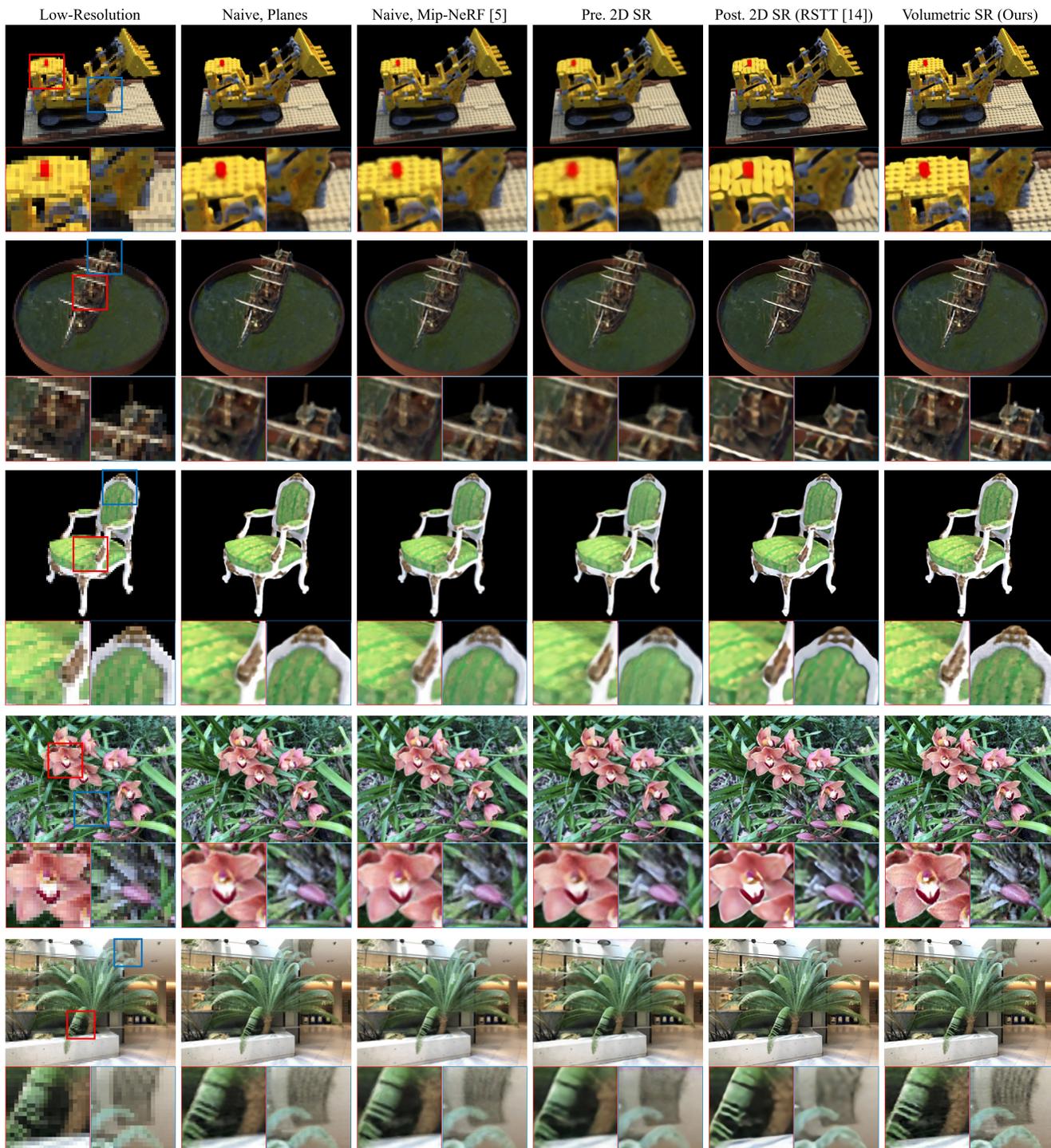


Figure 5. **Qualitative super-resolution evaluation.** We visually compare the performance of our method (right) to performance by leading SR methods and baselines, when super resolving novel low-resolution views (left). Our Neural Volume SR approach is able to reconstruct finer details and smaller structures compared to all other methods. While merely interpolating the feature planes in the volumetric domain (2nd from right) already yields a quality gain, applying our SR model F on the planes achieves a significant improvement.

the post-SR case, we also compare to using the prominent

SRGAN [22] method³ or the EDSR method [25], whose

³We use the distortion-minimizing training configuration (without an

Approach	Method	PSNR (dB)(\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
Naive	Mip-NeRF [5]	27.5	0.892	0.165
	Planes	27.4	0.892	0.153
pre. 2D SR	SwinIR [24]	25.6	0.872	0.183
post. 2D SR	EDSR [25]	25.9	0.876	0.148
	SRGAN [22]	24.5	0.863	0.149
	SWinIR [24]	25.1	0.863	0.183
	RSTT [14]	26.3	0.877	0.139
Volume. SR	Ours	28.5	0.911	0.099

Table 1. **Evaluation on synthetic scenes.** We evaluate all methods on 800 test images corresponding to novel views from 4 scenes unseen during training. The reported SSIM and PSNR (higher is better) and LPIPS (lower is better) values reveal a significant advantage of the proposed volumetric SR approach over all baseline approaches. Please refer to the text for details.

network architecture we employ for our model F . Finally, since one could view the LR scene captures as pertaining to an LR video sequence, we further experiment with post-SR using a state-of-the-art method for video spatio-temporal SR [14]. For the comparison to be unbiased, we train all 2D SR models from scratch using the same images from \mathcal{S}_t that we use for training our method.

Results. The quantitative comparisons reported in Tabs. 1 and 2 indicate the significant advantage of the proposed approach in terms of reconstruction quality (PSNR, SSIM) and perceptual quality (LPIPS) over all baselines, when evaluated on synthetic and real world scenes alike.

An additional inherent advantage of our approach, which super-resolves the volumetric representation itself, is its guaranteed across-view consistency. This advantage does not reflect in the quantitative results brought here, nor can it be conveyed in the still images depicted in Fig. 5. We therefore kindly refer readers to the videos provided in the Supp., which demonstrate this advantage over other baselines that perform SR in post-processing.

Interestingly, the results suggest that utilizing 2D SR methods (trained using the same data) actually has a negative effect on reconstruction performance, even compared to the naive baselines. We attribute this to the significant inconsistency resulting from independently processing different scene views.

Finally, we demonstrate the advantage of our approach visually in Fig. 5, where we compare novel test scene views rendered by our approach vs. the top performing baselines.

Ablations and Limitations. We further examine the effect of different components and hyper-parameter choices in our framework in a series of ablation experiments. We

adversarial loss) from SRGAN to train their SRResNet architecture, to allow a fair comparison with our method, which is trained to minimize distortion with respect to the GT images.

Approach	Method	PSNR (dB)(\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
Naive	Mip-NeRF [5]	21.9	0.657	0.401
	Planes	22.8	0.714	0.373
pre. 2D SR	SwinIR [24]	21.7	0.659	0.430
post. 2D SR	EDSR [25]	21.8	0.664	0.389
	SRGAN [22]	21.7	0.657	0.386
	SWinIR [24]	21.6	0.650	0.437
	RSTT [14]	22.0	0.679	0.377
Volume. SR	Ours	23.4	0.747	0.255

Table 2. **Evaluation on real-world scenes.** All methods are evaluated on all available 45 images corresponding to two real world scenes. For obtaining the ‘Naive: Planes’ and ‘Ours’ results, we use our model trained on views corresponding to other six scenes. Please refer to the text for details.

find that the resolution of the positional feature planes should be set to 1-2 times the resolution of the captured images set, while diverging from this range by a factor of 4 (in either direction) yields a drop of ~ 1.7 dB in PSNR, as we demonstrate visually in Fig. 4. Similarly, we analyse the significance of the 4th (view-direction) plane P_{dir} introduced in this work, and find that it is responsible to an increase of almost 3 dB in PSNR, which aligns with the visual demonstration in Fig. 3. Please refer to the Supp. for detailed results and additional ablation experiments. Finally, like most super-resolution methods, ours too is susceptible to non-ideal capturing conditions. We analyze their effect on performance compared to the other baselines in the Supplementary.

5. Conclusion

In this work, we introduce and investigate the novel task of neural volume super-resolution – rendering high-resolution views corresponding to a scene captured at low resolution. We propose a neural scene representation model that utilizes 2D feature planes, which we learn using LR scene captures. We then render novel, high-resolution views corresponding to the scene by super-resolving the 2D feature planes using a CNN model trained on a dataset of 3D scenes. We validate our approach both quantitatively and qualitatively on a diverse set of unseen synthetic and real-world scenes, and show a significant increase in performance compared to various baseline approaches in terms of view consistency, as well as in terms of reconstruction and perceptual quality, which demonstrate the major advantages of operating in the volumetric representation domain. Beyond the advantage in performing volumetric scene super-resolution, our work reveals the potential of this approach for reconstructing scenes captured under other types of realistic degradations (*e.g.* noise, blur), which we hope to investigate in future work.

Acknowledgment

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 945422.

References

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building rome in a day. *Communications of the ACM*, 54(10):105–112, 2011. [3](#)
- [2] Yuval Bahat, Netalee Efrat, and Michal Irani. Non-uniform blind deblurring by reblurring. In *Proceedings of the IEEE international conference on computer vision*, pages 3286–3294, 2017. [1](#), [2](#)
- [3] Yuval Bahat and Michal Irani. Blind dehazing using internal patch recurrence. In *2016 IEEE International Conference on Computational Photography (ICCP)*, pages 1–9. IEEE, 2016. [2](#)
- [4] Yuval Bahat and Tomer Michaeli. Explorable super resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2716–2725, 2020. [3](#)
- [5] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, October 2021. [3](#), [6](#), [8](#), [11](#), [12](#), [13](#)
- [6] 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](#)
- [7] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 417–424, 2000. [2](#)
- [8] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR’05)*, volume 2, pages 60–65. Ieee, 2005. [1](#), [2](#)
- [9] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. [2](#), [3](#), [4](#), [11](#)
- [10] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014. [2](#)
- [11] Runmin Dong, Lixian Zhang, and Haohuan Fu. Rrsrgan: Reference-based super-resolution for remote sensing image. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–17, 2021. [3](#)
- [12] Sina Farsiu, M Dirk Robinson, Michael Elad, and Peyman Milanfar. Fast and robust multiframe super resolution. *IEEE transactions on image processing*, 13(10):1327–1344, 2004. [2](#)
- [13] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5501–5510, June 2022. [3](#)
- [14] Zhicheng Geng, Luming Liang, Tianyu Ding, and Ilya Zharkov. Rstt: Real-time spatial temporal transformer for space-time video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17441–17451, 2022. [3](#), [8](#), [11](#), [13](#)
- [15] Daniel Glasner, Shai Bagon, and Michal Irani. Super-resolution from a single image. In *2009 IEEE 12th international conference on computer vision*, pages 349–356. IEEE, 2009. [2](#)
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [2](#)
- [17] Xin Huang, Qi Zhang, Ying Feng, Hongdong Li, Xuan Wang, and Qing Wang. Hdr-nerf: High dynamic range neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18398–18408, 2022. [3](#)
- [18] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. [3](#)
- [19] Jireh Jam, Connah Kendrick, Kevin Walker, Vincent Drouard, Jison Gee-Sern Hsu, and Moi Hoon Yap. A comprehensive review of past and present image inpainting methods. *Computer vision and image understanding*, 203:103147, 2021. [2](#)
- [20] Armin Kappeler, Seunghwan Yoo, Qiqin Dai, and Aggelos K Katsaggelos. Video super-resolution with convolutional neural networks. *IEEE transactions on computational imaging*, 2(2):109–122, 2016. [3](#)
- [21] Bahjat Kawar, Gregory Vaksman, and Michael Elad. Stochastic image denoising by sampling from the posterior distribution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1866–1875, 2021. [1](#), [2](#)
- [22] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017. [1](#), [3](#), [7](#), [8](#), [13](#)
- [23] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. Benchmarking single-image dehazing and beyond. *IEEE Transactions on Image Processing*, 28(1):492–505, 2018. [2](#)
- [24] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration us-

- ing swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1833–1844, 2021. [2](#), [6](#), [8](#), [13](#)
- [25] 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. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [13](#)
- [26] Ding Liu, Zhaowen Wang, Yuchen Fan, Xianming Liu, Zhangyang Wang, Shiyu Chang, and Thomas Huang. Robust video super-resolution with learned temporal dynamics. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2507–2515, 2017. [3](#)
- [27] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. [3](#)
- [28] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Srflo: Learning the super-resolution space with normalizing flow. In *European conference on computer vision*, pages 715–732. Springer, 2020. [3](#)
- [29] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. [3](#)
- [30] Tomer Michaeli and Michal Irani. Nonparametric blind super-resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 945–952, 2013. [1](#), [2](#)
- [31] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16190–16199, 2022. [3](#)
- [32] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. [1](#), [3](#), [4](#), [5](#), [6](#), [14](#), [15](#)
- [33] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2856–2865, 2021. [1](#)
- [34] Naama Pearl, Tali Treibitz, and Simon Korman. Nan: Noise-aware nerfs for burst-denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12672–12681, 2022. [3](#)
- [35] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *European conference on computer vision*, pages 501–518. Springer, 2016. [3](#)
- [36] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7495–7504, 2021. [1](#)
- [37] Jian Sun, Zongben Xu, and Heung-Yeung Shum. Image super-resolution using gradient profile prior. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. [2](#)
- [38] Yu-Wing Tai, Shuaicheng Liu, Michael S Brown, and Stephen Lin. Super resolution using edge prior and single image detail synthesis. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 2400–2407. IEEE, 2010. [2](#)
- [39] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. [3](#)
- [40] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Ji-aya Jia. Scale-recurrent network for deep image deblurring. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8174–8182, 2018. [1](#), [2](#)
- [41] 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. [3](#)
- [42] Bartłomiej Wronski, Ignacio Garcia-Dorado, Manfred Ernst, Damien Kelly, Michael Krainin, Chia-Kai Liang, Marc Levoy, and Peyman Milanfar. Handheld multi-frame super-resolution. *ACM Transactions on Graphics (TOG)*, 38(4):1–18, 2019. [2](#)
- [43] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. [6](#)
- [44] Haitian Zheng, Mengqi Ji, Haoqian Wang, Yebin Liu, and Lu Fang. Crossnet: An end-to-end reference-based super resolution network using cross-scale warping. In *Proceedings of the European conference on computer vision (ECCV)*, pages 88–104, 2018. [3](#)

Appendices

A. Additional Evaluations

This section provides additional qualitative and quantitative evaluations that evaluate the proposed method in comparison to existing methods for volumetric super-resolution.

A.1. Video Comparisons for Evaluating Across View-direction Consistency

The results in Tab. 1 demonstrate the advantage of our approach in terms of reconstruction quality, as measured independently on each rendered novel view. However, a key advantage of directly super-resolving the volumetric representation of the scene is that the added high-frequency, fine details are consistent across different views. While this consistency is hard to assess when examining each output independently, it becomes apparent when merging the outputs into a video trajectory, where geometric inconsistencies across adjacent viewing direction manifest themselves as prominent fidgeting-like temporal artifacts. Therefore, to demonstrate the advantage of our volumetric SR approach, we compile a series of ‘fly-over’ videos corresponding to the rendered scenes. Please refer to the attached video files ‘RealWorldScenes.mp4’ and ‘SyntheticScenes.mp4’ to view side-by-side comparisons of results rendered through different methods operating in 3D (including ours) with results by methods that apply 2D SR (including video SR [14]) in post processing. Please note the sharper results produced by our method, as well as the prominent jittering artifacts introduced by 2D SR post processing baselines.

A.2. Effect of Non-Ideal Capturing Conditions

Like most super-resolution methods, ours too is susceptible to degradations such as noise or non-ideal (typically unknown) downsampling kernels (PSFs). We perform experiments to quantify the effect of these conditions on the performance of our method, compared with the performance of the top performing baseline methods from Tab. 1 under the same settings.

To this end, we simulate capturing with noise (typically under low-lighting conditions) by adding random Gaussian noise with $\sigma = 5$ to the set of LR image captures \mathcal{I}_{lr}^s corresponding to a test scene s . Similarly, we simulate the non-ideal PSF case in a separate experiment by using a Gaussian blur kernel with a standard deviation of $\sigma = 1$ pixel for downsampling the GT HR images in \mathcal{I}_{hr}^s when creating the corresponding set of LR inputs \mathcal{I}_{lr}^s .

We report the results from the unknown PSF and additive noise experiments to the left and right of the slash sign in Tab. 3, respectively. Compared to the ideal setting case, the performance of all three compared baselines is reduced considerably, bringing them to perform at par under these conditions. In the case of our method, we conjecture that the reduced performance can mainly be attributed to the undermining effect that the non-ideal capturing may have over the minimization of the \mathcal{L}_{incon} loss term. We believe addressing these interesting challenges can be a topic for future work.

Approach	Method	PSNR (dB)(\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
Naive	Mip-NeRF [5]	26.68/26.87	<u>0.872/0.762</u>	0.205/ 0.190
	Planes	26.27/26.25	0.869/ 0.800	<u>0.201/0.260</u>
Volume. SR	Ours	<u>26.64/26.60</u>	0.875/0.726	0.189/0.249

Table 3. **Non-ideal capturing conditions for \mathcal{I}_{lr}^s .** Performance of our method compared with the leading baselines from Tab. 1 when simulating unknown, non-ideal downsampling kernel (left values) or additive noise (right values). We use a Gaussian blur downsampling kernel (with $\sigma = 1$ pixel) and an additive Gaussian noise (with $\sigma = 5$ gray levels) for creating the LR image sets in these simulations, respectively. Experiment settings are identical to those in Tab. 1. Please refer to the text for more details.

B. Ablation Experiments

To understand the significance of different elements and design choices in our framework, we perform a series of ablation studies. We next elaborate on the settings, report the results and derive conclusion from each of these experiments.

B.1. Significance of View-direction (4th) Plane

As we describe in Sec. 3.1, Chan *et al.* [9] use a tri-plane model constituting three positional feature planes $\{P_{xy}, P_{xz}, P_{yz}\}$ for representing human face scenes. We propose to add a 4th feature plane P_{dir} whose two axes span the 2D space of possible viewing directions for each point p . We do that in order to accommodate effects like non-uniform scene illumination, and subsequently to improve reconstruction performance, as we show in Fig. 3.

To further evaluate the effectiveness of P_{dir} , we compare the proposed quadri-plane model with a tri-plane model comprising only the three positional planes. Note that we did not include SR model F in this ablation experiment. We present a qualitative side-by-side comparison in the attached ‘EffectOfViewdirPlane.mp4’ video file, and provide a quantitative performance evaluation in Tab. 4, where we measure the performance of the two models on the task of reconstructing 1400 novel views corresponding to 7 scenes, using the PSNR, SSIM and LPIPS metrics. The results by all metrics, including an almost 3dB increase in PSNR, indicate a clear advantage for the proposed quadri-plane model.

P_{dir}	PSNR (db)(\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
With	28.7	0.898	0.078
Without	25.8	0.867	0.122

Table 4. **Effect of 4th feature plane P_{dir} .** Reconstruction performance of our representation model G with vs. without utilizing the 4th view-direction plane P_{dir} . Results indicate a significant advantage for our quadri-plane (top row) as features in P_{dir} are able to account for view-dependent effects such as non-uniform illumination and reflections.

B.2. Effect of Feature Planes Resolution N

The spatial dimensions N of the positional feature planes in the proposed quadri-plane representation model strongly affects representation performance, as we discuss in Sec. 3.1 and demonstrate in Fig. 4. Too small planes cannot represent high visual frequencies. On the other hand, using N values much larger than the resolution of the captured scene images results in visual artifacts when rendering novel views, as some plane regions remain uncovered by the information in \mathcal{I}^s . As a rule of thumb, we set N to be one to two times the resolution of the captured images $I \in \mathcal{I}^s$. To validate this choice, we compare reconstruction performance over 1400 novel views corresponding to 7 scenes when setting N to be a quarter, the same as or 4 times the resolution of I . The representations in this experiment were learned using images $I \in \mathcal{I}^s$ of size 400×400 , and we use each learned feature plane as is, without super-resolving it with F . The results in Tab. 5 suggest setting an appropriate resolution value N has a significant effect, and support our rule of thumb.

$N/\text{resolution}(I \in \mathcal{I}^s)$	PSNR (db)(\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
1/4	27.3	0.894	0.126
1	28.8	0.920	0.075
4	27.2	0.879	0.127

Table 5. **Effect of positional feature planes resolution N .** Performance of our representation model when using positional feature planes of different sizes N . Results suggest this parameter strongly effect performance, as supported by the qualitative comparison in Fig. 4.

B.3. Effects of Inconsistency Loss Term $\mathcal{L}_{\text{incon}}$ and SR Model F

We conduct additional ablation studies to analyse the role of these central components in the proposed volumetric SR approach. To this end, we repeat the evaluation experiment reported in Tab. 1 twice more: Once without performing the second training phase described in Sec. 3.3 (denoted ‘No $\mathcal{L}_{\text{incon}}$ ’ in Tab. 6), and a second time where we discard the feature-planes SR model F (denoted by ‘No SR’). In the latter experiment we perform the second training phase while optimizing only the weights of the decoder pair D , which take as input feature vectors v_p extracted directly from the learned low-resolution feature planes.

We report the results in Tab. 6, where we also bring as reference the results by the Mip-NeRF [5] baseline, by the naive baseline utilizing our quadri-plane model (denoted ‘LR feature planes (Naive)’ and the results by our full proposed configuration (bottom row). Note the significant drop in performance when eliminating either the SR model F (-0.7dB in PSNR) or the practice of minimizing $\mathcal{L}_{\text{incon}}$ (-1.4dB in PSNR), indicating the crucial role played by both of them.

C. Quantifying Across View-direction Consistency

To allow quantitatively evaluating this advantage of our approach, we next introduce a new Across View-direction Inconsistency (AVI) metric which aims to quantify the degree of geometrical inconsistency between every pair of test frames corresponding to adjacent viewing directions. To this end, we harness the GT optical flow that can be extracted for every adjacent frames pair in our synthetically rendered scenes.

Approach	Method	PSNR (dB)(\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
NeRF	Mip-NeRF [5]	27.5	0.892	0.165
Volumetric SR (Ours)	No SR	27.8	0.899	0.126
	No $\mathcal{L}_{\text{incon}}$	27.1	0.887	0.160
	LR Feature Planes (Naive)	27.4	0.892	0.153
	Proposed Configuration	28.5	0.911	0.099

Table 6. **Ablation experiments.** Comparing performance by the proposed method (bottom row) with its ablated variants when either excluding feature planes SR model F (denoted ‘No SR’), or not minimizing $\mathcal{L}_{\text{incon}}$. Experiment settings are identical to those in Tab. 1. The significant performance drop indicate the importance of both of components to the performance of our method. Please refer to the text for more details.

Approach	Method	AVI(\downarrow)
Naive	Mip-NeRF [5]	1.63
	Planes	1.97
pre. 2D SR	SwinIR [24]	1.71
post. 2D SR	EDSR [25]	2.82
	SRGAN [22]	2.92
	SWinIR [24]	2.61
	RSTT [14]	2.51
Volume. SR	Ours	1.86

Table 7. **Across View-direction Inconsistency (AVI).** Inconsistency scores (lower is better) are evaluated across 800 test frames corresponding to four synthetic scenes, following the experiment configuration from Tab. 1. The results reveal the inherent advantage of methods operating in the 3D domain (including ours) over methods performing 2D SR in post processing. Please refer to the text for more details.

We bring the formula for computing the AVI score and describe in detail each of its parts in Sec. C.1. We report AVI scores corresponding to our approach, as well as to all other compared baselines, in Tab. 7. The results reveal the inherent advantage of methods operating in the 3D domain. In terms of consistency, our method scores in between those baselines and significantly better than all methods performing 2D SR in post processing.

C.1. Calculating the AVI Score

We use the following formula to calculate the AVI score corresponding to a scene s :

$$\text{AVI} = \frac{1}{n} \sum_{k=0}^{\|\hat{\mathcal{I}}_{\text{sr}}\|-1} M_{k,k+1} \cdot \|\mathbf{g}(\text{Warp}(\hat{i}_{\text{sr}}^k)) - \mathbf{g}(\hat{i}_{\text{sr}}^{k+1})\|_2. \quad (6)$$

The AVI for each scene is computed over the set $\hat{\mathcal{I}}_{\text{sr}}$ of rendered HR views \hat{i}_{sr}^k corresponding to all test frames available for this scene. Since we have access to the Blender models used to generate the synthetic scene datasets, we are able to compute ground truth optical flow motion vectors both for the forward flow $f_{k \rightarrow k+1}$ and for the backward flow $f_{k+1 \rightarrow k}$ between all adjacent view-direction pairs $k, k+1$. We utilize the *backward flow* to warp each image to the next adjacent viewing direction via

$$\text{Warp}(\hat{i}_{\text{sr}}^k)(r) = \hat{i}_{\text{sr}}^k(r + f_{k+1 \rightarrow k}(r)), \quad (7)$$

where r denotes the 2D pixel location.

Ideally we would want the local differences between each frame and its warped adjacent view direction counterpart to be independent of whether a scene is smooth or highly textured. To this end, we apply a pre-processing operator \mathbf{g} on each of the two frames independently, before calculating the difference between them:

$$\mathbf{g}(i) = \text{PatchNorm}(\text{Im2patch}(i)). \quad (8)$$

Here $\text{Im2patch}(\cdot)$ extracts all overlapping 7×7 patches from an RGB image and $\text{PatchNorm}(\cdot)$ acts by discarding very smooth patches with standard deviation below a minimum value $\epsilon = 3/255$ (which are typically governed by noise) and



Figure 6. **Effect of Mask** $M_{k,k+1}$. We use masking to remove ghosting artifacts (e.g. the ship’s masts) in the warped image (a) prior to calculating the error. Here we illustrate the effect of using the mask with or without processing it with morphological image operators (b-d), and show that applying the fully processed mask yields significantly reduced ghosting artifacts (d).

normalizing the standard deviation of all other patches:

$$\text{PatchNorm}(p) = \begin{cases} 0, & \text{if } \text{std}(p) < \epsilon \\ p/\text{std}(p), & \text{otherwise} \end{cases}. \quad (9)$$

Even though we have access to the ground truth motion vectors, the warping can still introduce ghosting artifacts due to occlusions, i.e., pixels in image k that have no correspondence in image $k + 1$. To prevent errors stemming from the ghosting artifacts from dominating the AVI metric, we introduce an error mask $M_{k,k+1}$, computed as follows: First, by using the *forward flow*, we count how many pixels get mapped to the same pixel r :

$$\text{Count}_{k,k+1}(r) = \#\{q \mid \lfloor q + f_{k \rightarrow k+1}(q) \rfloor = r\}. \quad (10)$$

Pixels from viewing direction $k + 1$ for which this count is zero do not have a correspondence in viewing direction k . Thus, our mask is defined as

$$\tilde{M}_{k,k+1}(r) = \begin{cases} 0, & \text{if } \text{Count}_{k,k+1}(r) = 0 \\ 1, & \text{otherwise} \end{cases}. \quad (11)$$

Since the floor operation in the count computation can still lead to some ghosting due to subpixel occlusions, we apply a morphological closing followed by erosion to the mask, using a ”cross” structuring element of radius 1. The final mask in Eq. 6 is thus given by

$$M_{k,k+1} = \text{erosion}(\text{closing}(\tilde{M}_{k,k+1})). \quad (12)$$

When computing the AVI score, we only consider unmasked elements for the mean computation, i.e. $n = \#\{r \mid M_{k,k+1}(r) = 1\}$ in Eq. 6. The effects of the raw mask and the different morphological operators are visualized in Fig. 6. Note that we use the same motion vectors, and thus also the same masks, when computing the AVI score for each of the compared methods.

D. Synthetic Dataset 3D Scene Models

Our approach learns to super-resolve the volumetric representations by training on a set of training scenes \mathcal{S}_t . Being a deep-learning-based approach, it requires sufficient amount of training data to be able to generalize from the training scenes to unseen scenes. To meet this data requirement, we augment the synthetic scenes from [32] with 8 additional sequences generated from the following freely available Blender scenes

- <https://free3d.com/3d-model/bugatti-chiron-2017-model-31847.html>
- <https://free3d.com/3d-model/black-dragon-rigged-and-game-ready-92023.html>
- <https://free3d.com/3d-model/professional-scene-with-coca-cola-bottle-57999.html>
- <https://free3d.com/3d-model/gibson-es-335-816888.html>

- <https://free3d.com/3d-model/fuzzy-bear--10429.html>
- <https://free3d.com/3d-model/harley-davidson-low-rider-84874.html>
- <https://free3d.com/3d-model/holiday-beach-cartoon-scene-431138.html>
- <https://free3d.com/3d-model/donut-503129.html>
- <https://www.blendswap.com/blend/30606>
- <https://www.blendswap.com/blend/30328>
- <https://www.blendswap.com/blend/30072>
- <https://www.blendswap.com/blend/29594>
- <https://www.blendswap.com/blend/29435>
- <https://www.blendswap.com/blend/29080>
- <https://www.blendswap.com/blend/29508>
- <https://www.blendswap.com/blend/29461>
- <https://www.blendswap.com/blend/17636>
- <https://www.blendswap.com/blend/16763> (we removed the tender)
- <https://www.blendswap.com/blend/16155>
- <https://www.blendswap.com/blend/13078>
- <https://www.blendswap.com/blend/5364>

Training, validation and test sequences were generated using the Blender script provided by Mildenhall *et al.* [32].