

# SPARF: Large-Scale Learning of 3D Sparse Radiance Fields from Few Input Images

Abdullah Hamdi<sup>1</sup>

Bernard Ghanem<sup>1</sup>

Matthias Nießner<sup>2</sup>

<sup>1</sup> King Abdullah University of Science and Technology (KAUST)

<sup>2</sup> Technical University of Munich (TUM)

## Abstract

Recent advances in Neural Radiance Fields (NeRFs) treat the problem of novel view synthesis as Sparse Radiance Field (SRF) optimization using sparse voxels for efficient and fast rendering [14, 40]. In order to leverage machine learning and adoption of SRFs as a 3D representation, we present SPARF, a large-scale ShapeNet-based synthetic dataset for novel view synthesis consisting of  $\sim 17$  million images rendered from nearly 40,000 shapes at high resolution ( $400 \times 400$  pixels). The dataset is orders of magnitude larger than existing synthetic datasets for novel view synthesis and includes more than one million 3D-optimized radiance fields with multiple voxel resolutions. Furthermore, we propose a novel pipeline (*SuRFNet*) that learns to generate sparse voxel radiance fields from only few views. This is done by using the densely collected SPARF dataset and 3D sparse convolutions. *SuRFNet* employs partial SRFs from few/one images and a specialized SRF loss to learn to generate high-quality sparse voxel radiance fields that can be rendered from novel views. Our approach achieves state-of-the-art results in the task of unconstrained novel view synthesis based on few views on ShapeNet as compared to recent baselines. The SPARF dataset will be made public with the code and models on the project website [SPARF.com](http://SPARF.com).

## 1. Introduction

Although we observe the surrounding world only as a stream of 2D images, it is undeniably 3D. The goal of recovering this underlying 3D from 2D observations has been a longstanding goal of computer vision. The task of inverting the rendering process that creates the 2D projections we observe by trying to construct the 3D world is known as Vision as Inverse Graphics (VIG) [8, 22, 26, 67]. With the emergence of deep learning applications in computer graphics and the availability of 3D datasets, several approaches address the 3D generation task directly from 3D, without relying on appearance [1, 19, 27, 41, 46, 61]. However, recent developments in differentiable rendering have refueled the VIG direction,

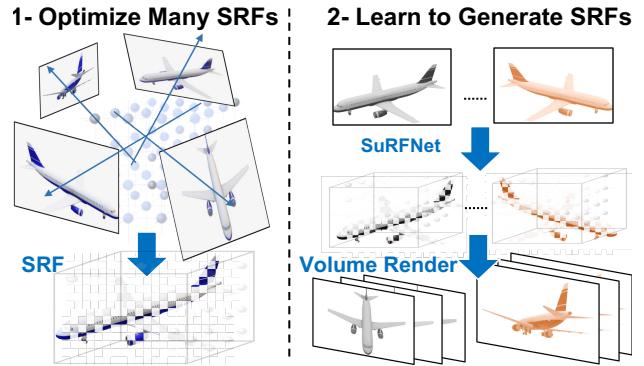


Figure 1. **Distribution of Radiance Fields.** We treat Sparse Radiance Fields (SRFs) as a 3D data structure, and learn the conditional generation of SRFs from few input images for the task of novel view synthesis. In order to do this, we build SPARF, a large-scale dataset of SRFs.

which facilitates using gradients of the rendering process to optimize for the underlying 3D setup based on image observations [15, 16, 23, 24, 28, 30, 31, 36, 38, 47, 53, 62, 68]. More specifically, Neural Radiance Fields (NeRFs) [7, 38, 45, 66] show impressive performance on novel view synthesis by optimizing volumetric radiance fields on a large number of posed multi-view images.

Various subsequent works addressed rendering speed [14, 40, 65], training size requirement [7, 39, 66], or pose requirements [34, 45] for NeRF. The seminal work of Plenoxyels showed that the MLP network is not necessary for quick optimization and volumetric rendering of the radiance fields. However, the current paradigm is still an optimization paradigm where a single scene representation is optimized without any generalization to new scenes/objects. In this work, we treat Sparse Radiance Fields (SRFs) as a 3D data structure and try to learn a generative model (dubbed *SuRFNet*) on the distribution of sparse-voxel radiance fields conditioned on a few images to generalize to unseen 3D shapes (see Figure 1).

In order to train deep learning models on 3D data to

Attribute	Posed Multi-View Datasets				
	SRN [51]	DTU [21]	NMR [43]	RTMV [52]	SPARF (ours)
Number of Classes	2	N/A	13	N/A	13
Number of Scenes/Objects	3,511	124	43,756	2,000	39,705
Image Resolution	128	1,200	64	1,600	400
Number of Radiance Fields	0	0	0	2,000	<b>1,072,008</b>
Real/Synthetic	synthetic	real	synthetic	synthetic	synthetic
View Setup	sphere	random	circle	hemisphere	sphere
Total Number of Images	265,550	4,235	1,050,144	300,000	<b>17,073,150</b>
Views per Model	50	N/A	24	N/A	430
Dataset Size (GB)	5.8	1	33	2,520	3,432

Table 1. **Comparison of Different Posed Multi-View Datasets.** We compare some of the widely used posed multi-view datasets to our large-scale SPARF dataset.

generalize to unseen examples, the dataset size should be in tens of thousands of samples [6, 58]. However, current posed multi-view datasets are not suitable for leveraging the power of deep networks as can be seen in Table 1. The image resolution is either too low (*e.g.*  $64 \times 64$  in [43]), the samples do not follow a controlled setup [21, 48], or lack diversity in the samples and classes [51]. For these reasons, we construct a large and high-resolution dataset (SPARF) of posed multi-view images from ShapeNet [6] that correspond to the same 13 classes originally used in the NMR dataset [43], but with an order of magnitude more images and pixels (17M *vs.* 1M images and  $400 \times 400$  *vs.*  $64 \times 64$  pixels). We also provide more than *one million* optimized sparse radiance fields of spherical harmonics and densities that allow for the novel view synthesis of the 40K models using Plenoxels [14].

The idea of learning a prior (2D CNN/ViT) on radiance fields in order to enhance the few-view setup of novel view synthesis is previously investigated by several works [29, 49, 66]. However, we propose SuRFNet to *directly* learn from the 3D sparse radiance fields, by optimizing partial SRFs from the few images and training a generalizable network that converts these partial SRFs to complete SRFs in a supervised fashion. Such a 3D setup benefits from structured 3D learning, creating a 3D prior that guarantees multi-view consistency, especially when rendering from out-of-distribution views. Also, this 3D sparse voxel setup benefits from the advancements in fast volume rendering [14, 40], allowing for end-to-end deep learning pipelines that harness volume rendering. To the best of our knowledge, our SURFNet is the first model that learns to generate 3D radiance fields for unseen objects at test time with only a few/single views by learning from the distribution of radiance fields in 3D.

**Contributions:** **(i)** To facilitate the application of deep learning on radiance fields, we provide a new Posed Multi-view dataset (SPARF) that is an order of magnitude larger than others (around 40K 3D models). The dataset includes a to-

tal of one million optimized Sparse Radiance Fields (SRFs) with multiple sparse voxel resolutions, which allows for high-quality novel view synthesis. **(ii)** We propose a novel architecture and a pipeline (SuRFNet) equipped with a specialized SRF-loss to generate voxel-based radiance fields from a few images based on learning to complete partial radiance fields. SuRFNet improves the performance of unconstrained novel view synthesis based on few views compared to state-of-the-art methods.

## 2. Related Work

**Learning 3D Shapes.** Several works aim to predict the geometry of 3D shapes given several input images, by directly optimizing the vertices of a template mesh through differentiable projections or through fitting a network [15, 16, 36, 53, 68]. Other works use MLPs as a deep prior to the optimized mesh [18, 36, 56]. On the other hand, some methods try to learn the distribution of 3D meshes by optimizing 3D generators independent of how the meshes look when rendered, solely based on the available 3D data and heuristic regularizers [10, 19, 41, 46]. Point cloud methods offer an alternative to the mesh complex topology by learning generative models on the point clouds themselves, *e.g.* by using an Auto Encoder [1, 61] or a GAN framework [1, 27]. The implicit representation paradigm offers an alternative to meshes for smooth and detailed shape representation. These methods learn a continuous implicit representation of shapes by learning the Signed Distance Functions or Occupancy of the object through MLPs [3–5, 30, 35, 42, 44, 50, 64]. While we learn a 3D representation in this work, the scope focuses on the quality of the rendering from novel views and not 3D reconstruction.  
**Neural Radiance Fields (NeRFs).** NeRFs [38] proved to be a successful popularizing in implicit volume representation and novel view synthesis. They define an implicit field and learn an MLP that predicts the RGB and density value of that 3D field given a set of posed images. NeRFs



Figure 2. **SPARF: a Large Dataset for 3D Shapes Radiance Fields and Novel Views Synthesis.**

shoot rays on the volume and integrate the predictions to obtain individual pixel values. This formulation, however, has many drawbacks including large memory and compute requirements, inability to model dynamic scenes, posed image requirements, and the limitation to small 3D objects or rooms [34, 42, 45, 65, 66]. To address the speed limitation, PlenOctreeNeRF [65] stores the precomputed RGB, density, and spherical harmonics in the 3D volume as an Octree data structure for fast inference. Plenoxels [14] optimize the density and spherical harmonics on sparse voxels with a TV loss and perform ray marching for rendering from novel views. Similarly, INGP [40] uses multi-resolution voxel hashing to perform a real-time rendering of radiance fields, demonstrating that the redundancy of the MLP in NeRFs. We build on these observations and build the SPARF dataset of sparse voxel radiance fields in order to facilitate learning on these SRFS as a 3D data structure instead of just a side outcome of a volumetric optimization.

**Few-Image NeRFs.** To address the original NeRF’s requirement of many posed images, several methods were proposed. The seminal work PixelNerf [66] is the first to reduce the image data requirements in order to learn a NeRF by using a trained CNN prior that can allow for transferable representation between scenes. Similarly, MVSNet [7, 17], AutoRF [39], and ShaRF [49] learn a CNN prior to generalize across scenes. IBRNet [54] learns to render novel views based on neighboring views and optimized neural volume representation. More recently, VisionNerf [29] proposes to use a ViT [12] to extract global features from the input images to enhance the capability of the NeRF MLP to predict the radiance field when one image is used as input. Unlike these works, we propose SuRFNet to directly learn from the 3D sparse radiance fields. Such a 3D setup benefits from structured 3D learning, creating a 3D prior that guarantees multi-view consistency while benefiting from the speed of recent voxel-based methods. A concurrent work by Guo *et al.* [17] learns a 3D prior based on a perceptual loss, but does

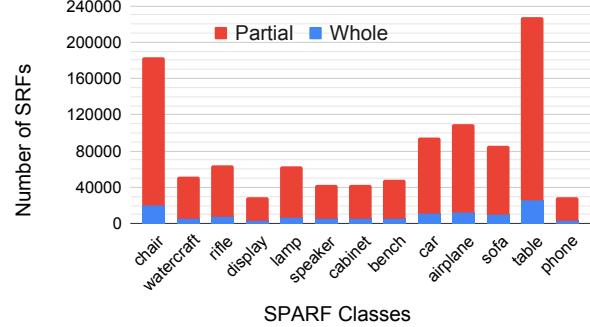


Figure 3. **SPARF Distribution.** We show the distribution of classes in SPARF and how the one million partial and whole SRFS are distributed. The numbers are equally distributed on three voxel resolutions: 512, 128, and 32.

not use 3D supervision and only uses dense voxels ( limiting the pipeline to the low resolution of  $64^3$  ).

**Datasets for novel view synthesis.** Several datasets were proposed to support the task of novel view synthesis. NeRF [38] introduced 8 synthetic scenes with 360-degree views. The rapid progress that followed this work was attributed to algorithmic developments rather than scaling up the data. However, for the successful application of deep learning, a large amount of data is necessary for improved generalization. Wang *et al.* [55] introduced a synthetic dataset including Google Scanned Objects. Other datasets for training multi-view algorithms include DTU [21], LLFF [37], Tanks and Temples [25], Spaces [13], RealEstate10K [70], SRN [51], Transparent Objects [20], ROBI [60], CO3D [48], SAPIEN [59], and BlendedMVS [63]. Recently, RTMV [52] introduced a ray-traced posed multi-view dataset with 2000 scenes and high-resolution images. Unfortunately, the current posed multi-view datasets commonly used in NeRF research are either small in image resolution or the number of posed images (SRN [51] and NMR [43]), small in the number of scenes/shapes and classes (Synthetic NeRFs [38]), or lack structure (DTU [21] and RTMV [52]). A detailed multi-attribute comparison is provided in Table 1.

### 3. SPARF: a Large Dataset of 3D Shape Radiance Fields (SRFs)

One of the goals of this work is to learn to generate high-quality Sparse Radiance Fields (SRFs) in one forward pass of a deep network to enable fast novel view synthesis. In order to do this, harnessing the power of deep networks would require a large dataset of SRFs in a controlled setup. We describe the details of SPARF in the following section.

#### 3.1. Dense Posed Multi-view Image Dataset

The first step in collecting the desired large high-quality radiance field dataset is to collect a synthetic posed multi-



Figure 4. **SRFs: The optimized Sparse Radiance Fields in SPARF.** A total of one million SRFs have been collected in SPARF, including on multiple voxel resolutions: 32 (top), 128 (middle), and 512 (bottom).

view dataset. We use ShapeNet Core 55 [6] as the data of choice for 3D shapes. For rendering, we used Pyglet [2] API through Trimesh library [11]. The renderer is based on OpenGL [57] rasterizer to render over 17 million images of around 40,000 shapes from 13 different classes at a high resolution of  $400 \times 400$ . Every shape is rendered equidistantly from 400 views distributed in a spherical configuration surrounding the object, including from the bottom (see Figure 2 for an example). An additional 20 views are rendered from random views from the same distance as test views for novel view synthesis tasks. Furthermore, an additional 10 views are rendered randomly from random distances bounded by a reasonable range, such that at least a part of the object is guaranteed to be visible. This last set is aimed at robustness purposes to test whether novel view synthesis methods can generalize to out-of-distribution posed views. More details about the rendering setup, including lighting and materials are provided in supplementary material.

### 3.2. Multi-Resolution 3D SRFs

Sparse Radiance Field (SRF) can be defined as a voxel grid of dimension  $1 + d$ , where  $d$  is the dimension of radiance colors  $\rho_{i,j,k} \in \mathbb{R}^d$  at that specific  $(i, j, k)$  indexed voxel in addition to one dimension for density  $\alpha_{i,j,k} \in \mathbb{R}$ . We assume that the grid is of size  $H$  in each of the three dimensions:  $\mathcal{X} \in \mathbb{R}^{H^3 \times (1+d)}$ . Since the SRF is sparse, it can be represented with the COO format [9] as a set of  $M$  tuples of positive integer coordinates  $\mathbf{c} \in \mathbb{Z}^+$  and features  $\mathbf{f} \in \mathbb{R}^{d+1}$  with the sparsity of  $1 - \frac{M}{H^3}$  as follows:

$$\mathcal{X}_{\text{non-empty}} = \{(\mathbf{c}_m, \mathbf{f}_m)\}_{m=1}^M \quad (1)$$

The ordering of the set of tuples is arbitrary, but the features  $\mathbf{f}_m$  consist of the density  $\alpha_{i,j,k}$  and radiance colors  $\rho_{i,j,k}$  at that location  $\mathbf{c}_m = (i, j, k)$ . For the radiance field colors, we use the spherical harmonics proposed in Plenoxels [14] for view-dependent learning of radiance common in NeRFs [38]. The SRF can be viewed as a distillation of the NeRF MLP into sparse voxels for efficient optimization and volume rendering. In many 3D object tasks, a coarse-to-fine approach is followed, demanding multiple resolutions. Hence, we collect the SPARF with multiple resolutions  $H \in \{32, 128, 512\}$ , as shown in Figure 4. We used an adaptation of Plenoxels [14] to collect the dataset of a total of one million SRFs as we detail next. In order to scale up the Plenoxels optimization for this huge number of shapes and variants, we utilize a large number of images in the SPARF dataset to reduce the iterations to a minimum number while maintaining a high average PSNR across the dataset for the collected SRFs across the multiple resolutions. A total of 200K GPU hours are used in the optimization process to collect SPARF. highly detailed 3D meshes can be extracted easily from the collected SRFs as can be seen in Figure 5.

### 3.3. Representing Images with Partial 3D Radiance Fields

In addition to collecting the “whole” part of SPARF that utilized all 400 images for every shape in optimizing the SRFs, we collect “partial” SRFs. These are SRFs optimized on only a small number of images (1 or 3) randomly sampled from all 400 images, resulting in multiple partial SRF variants of that shape (see Figure 6). However, the same 3D



Figure 5. **Extracting 3D Meshes from SRFs.** Since SPARF and SuRFNet live on the 3D voxel’s space, extracting the mesh is straightforward with one pass of MarchingCubes [32].

object can have multiple partial SRFs depending on the input images optimized. Therefore, we collect multiple variants of this partial SRFs for every object. We use these *partial* SRFs as input to our training pipeline.

#### 4. Learning to Generate SRFs

Unlike previous methods that try to embed priors in learning to generate radiance fields from a few images (*e.g.* PixelNerf [66] and VisionNerf [29]), we distill the views into a 3D SRF and then perform the learning in the 3D sparse voxel space. SuRFNet offers an alternative and a new way to learn novel view synthesis by learning to generate the entire 3D radiance field based on a small observation of the scene. Such a 3D setup benefits from structured 3D learning, thus creating a 3D prior that guarantees multi-view consistency, especially when rendering from out-of-distribution views (as we show in Section 5.3). The input to the pipeline is the input partial SRFs from Section 3.3, where the goal is learning a generalizable network that converts partial SRFs to whole SRFs as can be seen in Figure 7.

##### 4.1. SuRFNet: Sparse Radiance Fields Network

**3D sparse conv.** We leverage the Minkowski Net [9] as the 3D sparse convolution network of choice. However, this type of sparse convolution is not designed for generative tasks of fine-grained details, which are the scope of this work. In a typical sparse convolution learning paradigm, the output itself is not a sparse voxel grid, but rather a point cloud or a continuous 3D prediction. In our setup, the output itself is a sparse voxel grid of radiance fields that have irregular structures and low-density components that cannot be seen when they are volume rendered yet they affect the SRF learning. We use residual connections in a U-Net form of Minkowski Net  $\mathbf{F}$  with  $l$  modules, where each module consists of 3 sparse convolutions layers with strides  $s$ .

**Challenges of Learning SRFs.** Even though the setup appears to be a simple encoder-decoder fully supervised learning setup from partial SRF to whole SRF, learning SRFs is much more challenging in reality. As a data structure, SRF is an irregular volumetric representation that does not necessarily reflect the underlying 3D shape/scene, as it is a result of the optimization of posed images into volume. Many of the non-empty voxels have low density and do not affect

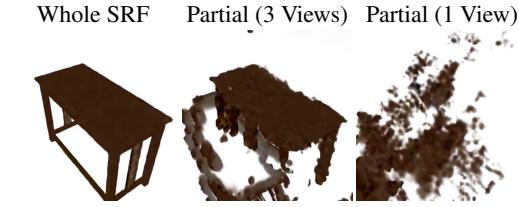


Figure 6. **Whole vs. Partial SRFs.** The partial SRFs are used instead of the few images that generated them as input to the learning pipeline to generate the whole SRFs

the volume rendering, but include color information that can confuse the network. Also, small errors in predicting the densities or radiance colors can result in large distortions in the rendered images, hurting overall novel view synthesis performance. Furthermore, the nature of the SRFs is closer to being a surface representation (usually the densities are low inside the object), which makes a useful signal to create the shape extremely sparse in the high-resolution 3D volume space. Another problem with sparse convolution is that a vanishing gradient is more imminent than in typical learning paradigms in 3D. The usual sparsity in our setup is  $\sim 99\%$ , and misalignment between the input SRF coordinates and output SRF coordinates can further harm the gradients and affect the learning process. In order to tackle these issues, we propose three specialized losses detailed next.

##### 4.2. SRF-Loss

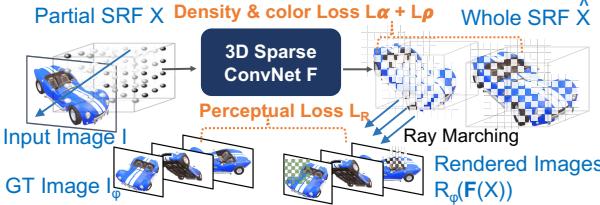
**Density loss.** The goal of the density loss is to create a dense surface. We propose the following binary cross-entropy loss on the predicted densities  $\alpha$  as follows:

$$L_\alpha(\mathcal{X}, \hat{\mathcal{X}}) = -(\hat{\mathbf{y}} \log(\mathbf{y}) + (1 - \hat{\mathbf{y}}) \log(1 - \mathbf{y})) \quad (2)$$

$$\text{s. t. } \hat{\mathbf{y}} = \mathbb{1}(\mathcal{S}(\hat{\mathcal{X}}_\alpha) > \alpha_{\text{dense}}), \mathbf{y} = \mathcal{S}(\mathbf{F}(\mathcal{X}))_\alpha$$

where  $\mathcal{X}, \hat{\mathcal{X}}$  are the input and output SRFs respectively (as defined in Eq (1)), and  $\alpha_{\text{dense}}$  is the density threshold distinguishing dense voxels from the air (usually set to 0). The sampling function  $\mathcal{S}$  samples points in the grid space where the loss on the outputs  $\mathbf{y}$  is defined. Defining the loss only on the non-empty voxels will leave many radiance clouds in the output, deteriorating the quality of the rendered images from the output SRF.

**Q-Gaussian loss sampling.** One of the challenges in working with voxels of high resolution (*e.g.* 512) is that the learning of the operations in sparse voxels cannot involve densifying the voxels to the original resolution (*e.g.*  $512^3$ ), due to prohibitive memory requirements. This is why the sampling function  $\mathcal{S}$  in Eq (2) is of utmost importance in guiding the training of SuRFNet. We sample at random coordinates centered at the center of the voxel grid  $\mathcal{S} : \mathbf{c} \sim \mathcal{Q}\left(\mathcal{N}\left(\frac{\mathbf{H}}{2}, \frac{H\sigma^2}{2}\mathbb{I}\right)\right)$ , where  $\mathbf{H} = (H, H, H)$  is the



**Figure 7. SuRFNet: Learning to Generate Whole Radiance Fields from Partial Views.** We process the input images into partial SRFs  $\mathcal{X}$  before learning a sparse convolutional network to generate the whole SRF. A perceptual loss is employed on the rendered images from poses  $\phi$  to enhance the perceptual quality of the generated SRF. The whole SRF  $\hat{\mathcal{X}}$  is used to 3D-supervise the SRF generation with density and radiance color losses.

voxel grid resolution vector,  $\mathbb{I}$  is the identity matrix,  $\sigma$  is a hyperparameter determining the spread of the loss, and  $\mathcal{Q} : \mathbb{R}^3 \rightarrow \mathbb{Z}^{+3}$  is the quantization-and-cropping function of coordinates that ensure the output coordinates are integers within bounds  $c \in [0, 1, \dots, H - 1]^3$ . We discuss more details about  $S$  and alternative configurations in Section 5.2. **Radiance color loss.** To ensure the output SRFs follow the ground truth optimized SRFs in radiance color, we follow the simple L1 loss on the radiance colors  $\rho$ . However, as mentioned earlier, some of the non-empty voxels contain low density and will not be seen in the rendering and can contain any random colors. Therefore, we mask these non-empty low-density voxels out of the L1 loss as follows:

$$L_\rho(\mathcal{X}, \hat{\mathcal{X}}) = \|\mathbf{M}_\alpha \mathbf{F}(\mathcal{X})_\rho - \mathbf{M}_\alpha \hat{\mathcal{X}}_\rho\|_1 \quad (3)$$

s. t.  $\mathbf{M}_\alpha = \mathbb{1}(\hat{\mathcal{X}}_\alpha > \alpha_{\text{dense}})$

**Perceptual loss.** Using only the 3D radiance color loss in Eq (3) ignores the rendering quality of the generated SRF, and would make it sensitive to hyperparameters (see Figure 11). Hence, we introduce an online perceptual loss that would volume render the generated SRF during training from  $M$  random views that come from the same ground truth image poses  $\phi$ , and an L1 loss is defined between the generated images and the ground truth posed images  $\mathbf{I}_\phi$ .

$$L_R(\mathcal{X}) = \|\mathcal{R}_\phi(\mathbf{F}(\mathcal{X})) - \mathbf{I}_\phi\|_1, \quad (4)$$

where  $\mathcal{R}_\phi$  is the fast volume rendering function that renders SRFs from poses  $\phi$  using the trilinear interpolation between voxels proposed in Plenoxels [14].

The final loss to train the network  $\mathbf{F}$  would be combining the three losses in Eq (2,3,4) as follows:

$$\text{Loss}_{\mathbf{F}} = L_\alpha + \lambda_\rho L_\rho + \lambda_R L_R, \quad (5)$$

where  $\lambda_\rho, \lambda_R$  are hyperparameters to control the radiance colors compared to the density predictions. The network is trained on all  $N$  whole SRFs  $\hat{\mathcal{X}}$  in the dataset, while the input SRFs  $\mathcal{X}$  are randomly chosen from several partial SRFs created by the same number of images from those shapes.



**Figure 8. SPARF vs. other Datasets.** SPARF offers a large-scale high-resolution dataset compared to other posed multi-view datasets. We show the same chair here on SPARF, SRN, and NMR (please zoom-in for differences). This highlights the huge quality gap between SPARF and other ShapeNet-based datasets.

## 5. Experiments

### 5.1. Collecting SPARF

The engineering aspect of collecting, storing, and organizing the one million SRFs with multiple resolutions is as challenging as training properly on SRFs. In order to do that in manageable time and memory, while maintaining high quality in the optimized samples, a set of strategies is employed. The dimension of the radiance color  $d$  is chosen to be  $d = 3 \times 4 = 12$  of 4 spherical harmonics factors of RGB channels for view-dependent SRF and  $d = 3 \times 1 = 3$  for fixed RGB colors of the SRF. Since the input partial SRFs are noisy, we use  $d = 3$  while the final output SRFs use  $d = 12$  for high-quality image generation. Using fewer Spherical Harmonics components (from 9 to 4 per RGB channel) reduces the optimization space by 40% and time by 10%, while maintaining the same PSNR. Using RGB as colors instead of SH factors reduces PSNR by  $\sim 1$  dB, space by 80%, and time by 20%. Running Plenoxels [14] for fewer iterations ( $3 \times 12K$ ) reduces the time by 30% while maintaining the same PSNR. Using 400 views/shapes in SPARF to optimize the SRFs keep the time manageable in optimization ( $\sim 4$  minutes for the 512 resolution) while maintaining high PSNR ( $\sim 30$  dB). When saving the SRFs, we only save the set of coordinates (integers) and float features (densities and radiance components). A total of four variants of the partial SRFs are collected for all the resolutions and the partials use 1 and 3 images. The anatomy of the distribution of classes and SRFs in SPARF is presented in Figure 3. More details about SPARF and visualizations of some of its samples are available in the supplementary material.

### 5.2. Training Setup

**Dataset.** We pick our SPARF for the task of predicting whole SRFs for the purpose of novel view synthesis. The other datasets (SRN [51] and NMR [43]) are too small or low in resolution, which prevents optimizing high-quality radiance fields (see Figure 8).

**Evaluation metrics.** Following the previous novel view synthesis works [14, 29, 66], we use PSNR, SSIM, and LPIPS

Baselines	SPARF Classes													mean
	chair	watercraft	rifle	display	lamp	speaker	cabinet	bench	car	airplane	sofa	table	phone	
Plenoxels [14] (1V)	9.2	11.1	11.7	8.0	13.6	8.2	10.4	10.5	7.1	12.8	9.3	9.9	8.3	10.0
Plenoxels [14] (3V)	10.7	13.3	14.9	9.7	15.8	10.4	12.4	11.6	7.1	14.6	11.6	10.8	9.7	11.7
PixelNerf [66] (1V)	13.3	16.3	16.7	11.9	17.6	11.3	14.5	14.6	13.2	19.2	13.3	13.2	13.2	14.5
PixelNerf [66] (3V)	13.5	16.6	16.9	12.2	17.9	11.9	14.9	14.8	13.4	19.4	13.4	13.3	13.3	14.7
VisionNeRF [29] (1V)	13.0	15.6	15.8	11.7	16.7	11.2	14.0	14.3	12.7	17.8	13.3	13.0	12.6	14.0
<b>SuRFNet (ours) (1V)</b>	11.6	16.2	17.0	12.0	16.2	12.6	17.0	13.5	16.6	17.5	14.1	10.1	15.3	14.6
<b>SuRFNet (ours) (3V)</b>	<b>15.3</b>	<b>18.3</b>	<b>18.8</b>	<b>15.0</b>	<b>19.0</b>	<b>16.6</b>	<b>20.0</b>	<b>15.6</b>	<b>16.6</b>	<b>18.5</b>	<b>18.1</b>	<b>14.9</b>	<b>17.8</b>	<b>17.3</b>

Table 2. **SPARF Benchmark on Out-of-distribution View Synthesis.** We compare the validation PSNR of some of the widely used novel view synthesis techniques on the SPARF dataset for the generalization of novel view synthesis beyond a single example and on view tracks completely different from the ones seen in training views. One view (1V) and three views (3V) inputs are reported.

[69] as metrics to evaluate the synthesis. However, one key difference between our work and previous ones is that our setup is a learning setup (with training and validation), while previous works treat it as an optimization problem. Most previous works on novel view synthesis try to only generalize the generated views on the *same shape*, while we aim to generalize *across shapes* and *across views*. We treat the collected whole SRFs as ground truth labels for the input few images from the training set. We consider the validation PSNR, SSIM, and LPIPS of the input images at the validation SRF set of shapes (on the test images of those shapes) as the main evaluation metrics. Also, we report validation accuracy =  $\frac{\text{validation PSNR with test few images}}{\text{whole SRF optimization's PSNR}}$  and propose it as a new metric to evaluate such a learning setup of SRFs. Furthermore, as we describe in Section 3, SPARF has 10 posed Out-Of-Distribution (OOD) images for every 3D shape to evaluate the robustness of novel view synthesis methods in the unconstrained setup. We report these *OOD* PSNR, SSIM, LPIPS, and Accuracy as well.

**Baselines.** We use PixelNerf [66], Plenoxels [14], and VisionNerf [29] as the main baselines for our work. Our SuRFNet network has two sizes: large (87 million parameters) and small (13 million parameters). More details can be found in the supplementary material.

### 5.3. Results

We show qualitative results of generating novel views from single input images on unseen shapes in Figure 9. We also show qualitative comparisons in Figure 10. The generated SRFs can be found in the supplementary material. We present a summary of the quantitative evaluations next, where SuRFNet achieves state-of-the-art results on unconstrained novel view synthesis from one or few images on unseen shapes.

**SPARF View-Generalization benchmark.** In Table 2, we report the average PSNR results on the validation set of SPARF for different methods on unseen shapes during training on all 13 different object classes and on out-of-distribution views. It shows that our SuRFNet can generalize to out-of-distribution views on unseen shapes during test



Figure 9. **SuRFNet: Generating High-Resolution Radiance Fields.** We show some volume-rendered sequences based on our SuRFNet voxel radiance field outputs (512 resolution), given only 3 images of each shape.

time, surpassing state-of-the-art PixelNerf [66] and VisionNerf [29]. Visual comparisons can be found in Figure 10. As can be seen from those results, the learned 3D prior results in multi-view consistency, especially when rendering from out-of-distribution views. More results can be found in supplementary material.

## 6. Analysis and Insights

### 6.1. Ablation Study

We ablate different components of SuRFNet’s architectures and the loss configuration choices and report the results in Tables 3 and 4. More ablations on the network, loss, and hyperparameters of training SuRFNet can be found in supplementary material.

**Training SuRFNet.** Results show that increasing the size of the network (from 13M to 87M parameters) helps improve generalization accuracy. Also, they show the importance of the loss components proposed in Eq (5). The use of only density loss creates a reasonable dense shape but without colors. While combining the density loss with the 3D radiance color loss creates colorful objects, it does not perform

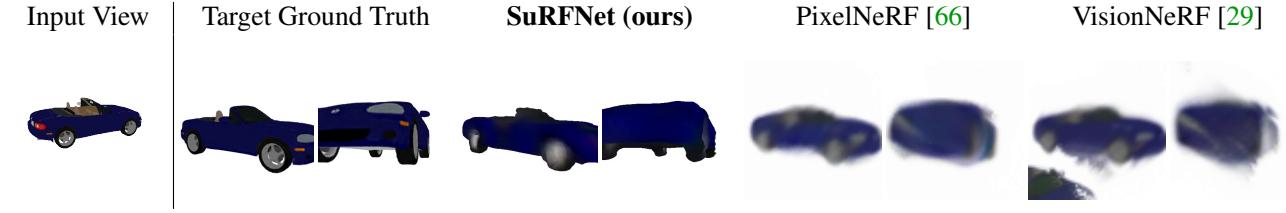


Figure 10. **Qualitative Comparisons.** We show different render from our SuRFNet outputs generated from a single image compared to other methods (pixel-Nerf [66], and VisionNerf [29]) and whole SRF "GT" renderings. Note that the predicted two views lay outside the training views distribution (zoomed in randomly). This test highlights the weakness of the 2D-based baselines [29, 66] outside the training track, while our 3D approach maintains multi-view consistency everywhere.

3D Backbone		Loss components			Results
Small	Large	$(L_\alpha)$	$(L_\rho)$	$(L_R)$	Val. Acc.
✓	-	✓	✓	-	65.2
✓	-	✓	✓	✓	65.7
-	✓	✓	✓	-	66.4
-	✓	✓	✓	✓	<b>68.2</b>

Table 3. **Ablation Study.** We ablate different components of in SuRFNet (3D backbone and SRF-Loss) and report validation accuracy of car class.

Strategy	1-view			3-view		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
uniform	20.03	0.94	0.10	21.00	0.94	0.09
Q-Gaus.	20.55	0.93	0.09	21.83	0.94	0.08

Table 4. **Effect of Loss Sampling Strategy.** We study the effect of Loss sampling strategy (uniform vs. Q-Gaussian) on airplane class.

well in the task of novel view synthesis as it does not respect how the object looks, and any deviations from the labeled radiance colors can cause large image distortions. Please see Figure 11 for the importance of the perceptual loss.

**Loss Sampling.** We study the effect of the sampling strategy with a different number of input images at test time on the performance of SuRFNet in Table 4. It shows that using a uniform sampling strategy depletes the learning capacity of the network and can degrade performance. The effect is more evident when the number of views is one, where the partial SRFs are more sparse and the training is delicate.

## 6.2. Speed and Compute Cost

To assess the contributions of the SuRFNet pipeline, we study the time and memory requirements of each element in the pipeline. We record in Table 5 the number of floating-point operations (GFLOPs) and the runtime of a forward pass (including rendering) for a single output image from one input image on Titan RTX GPU.

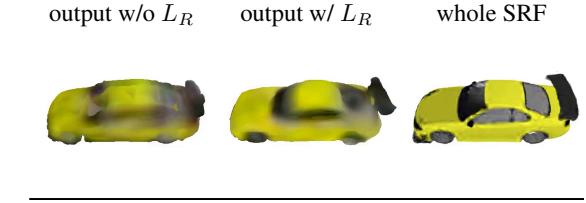


Figure 11. **Effect of the Perceptual Loss  $L_R$ .** Adding a perceptual loss on volume-rendered images during training SuRFNet insures the rendered images remain closer to how they should be rendered, as the 3D radiance colors supervision won't guarantee the rendering quality. (*left*): without perceptual loss , (*middle*): with the loss.

Network	Network FLOPs (G)	Network Inference (ms)	Parameters Number (M)	Rendering Speed (FPS)
PixelNeRF [66]	7.3	5.33	21.8	1.2
VisionNerf [29]	33.7	12.5	68.6	1.2
SuRFNet (small)	~15	14.4	13.4	15
SuRFNet (large)	~100	90.0	87.3	15

Table 5. **Time and Memory Requirements.** We assess the computational cost of the main components studied

## 7. Conclusions and Future Work

We propose a large-scale dataset SPARF of sparse radiance fields that include around one million SRFs and 17 million posed images of 3D shapes. The dataset aims to move the community in the direction of treating radiance fields as a 3D data structure, instead of optimization results and MLP fitting. Leveraging the utility of SPARF, we propose a SuRFNet pipeline to train a conditional generative model to generate SRFs from few input images (1 or 3) distilled as partial SRFs. SuRFNet allows generating radiance fields from single images of unseen shapes, which allows for rendering high-quality images from novel views, reaching state-of-the-art performance in unconstrained novel view synthesis compared to other methods.

One crucial limitation in this work is the large amount of compute and memory necessary to create, store, and process SRFs, especially at high-resolution voxel grids. This creates a bottleneck in training, developing, and building

on SRFs. Developing efficient methods to work and learn from sparse voxel grids would be a viable plan moving forward in order to develop deep and large models in this space, as well as, leveraging popular (slow) generative models (*e.g.* diffusion models) on 3D radiance fields.

## References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. *International Conference on Machine Learning (ICML)*, 2018. [1](#), [2](#)
- [2] Alex Holkner et al. Pyglet. [4](#)
- [3] Shengqu Cai, Anton Obukhov, Dengxin Dai, and Luc Van Gool. Pix2nerf: Unsupervised conditional p-gan for single image to neural radiance fields translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3981–3990, June 2022. [2](#)
- [4] 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, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16123–16133, June 2022. [2](#)
- [5] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5799–5809, 2021. [2](#)
- [6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. [2](#), [4](#), [14](#)
- [7] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. [1](#), [3](#)
- [8] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. In *Advances in Neural Information Processing Systems*, pages 9609–9619, 2019. [1](#)
- [9] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. [4](#), [5](#), [12](#)
- [10] Angela Dai and Matthias Nießner. Scan2mesh: From unstructured range scans to 3d meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5574–5583, 2019. [2](#)
- [11] Dawson-Haggerty et al. trimesh. [4](#)
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. [3](#)
- [13] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2367–2376, 2019. [3](#)
- [14] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinzhong 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. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [12](#), [14](#)
- [15] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9785–9795, 2019. [1](#), [2](#)
- [16] Shubham Goel, Georgia Gkioxari, and Jitendra Malik. Differentiable stereopsis: Meshes from multiple views using differentiable rendering. *arXiv preprint arXiv:2110.05472*, 2021. [1](#), [2](#)
- [17] Pengsheng Guo, Miguel Angel Bautista, Alex Colburn, Liang Yang, Daniel Ulbricht, Joshua M Susskind, and Qi Shan. Fast and explicit neural view synthesis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3791–3800, 2022. [3](#)
- [18] Rana Hanocka, Gal Metzger, Raja Giryes, and Daniel Cohen-Or. Point2mesh: A self-prior for deformable meshes. *arXiv preprint arXiv:2005.11084*, 2020. [2](#)
- [19] Qixing Huang, Xiangru Huang, Bo Sun, Zaiwei Zhang, Junfeng Jiang, and Chandrajit Bajaj. Arapreg: An as-rigid-as possible regularization loss for learning deformable shape generators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5815–5825, 2021. [1](#), [2](#)
- [20] Jeffrey Ichniowski, Yahav Avigal, Justin Kerr, and Ken Goldberg. Dex-nerf: Using a neural radiance field to grasp transparent objects. *arXiv preprint arXiv:2110.14217*, 2021. [3](#)
- [21] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014. [2](#), [3](#)
- [22] Danilo Jimenez Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4996–5004. Curran Associates, Inc., 2016. [1](#)
- [23] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–386, 2018. [1](#)
- [24] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of the IEEE Conference*

- on Computer Vision and Pattern Recognition (CVPR)*, pages 3907–3916, 2018. 1
- [25] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 3
- [26] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *Advances in neural information processing systems (NIPS)*, pages 2539–2547, 2015. 1
- [27] Ruihui Li, Xianzhi Li, Ka-Hei Hui, and Chi-Wing Fu. Sp-gan: Sphere-guided 3d shape generation and manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021. 1, 2
- [28] Tzu-Mao Li, Miika Aittala, Frédéric Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing through edge sampling. In *SIGGRAPH Asia 2018 Technical Papers*, page 222. ACM, 2018. 1
- [29] Kai-En Lin, Lin Yen-Chen, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. Vision transformer for nerf-based view synthesis from a single input image. In *WACV*, 2023. 2, 3, 5, 6, 7, 8, 14, 21, 22
- [30] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019. 1, 2
- [31] Matthew M Loper and Michael J Black. Opendr: An approximate differentiable renderer. In *European Conference on Computer Vision (ECCV)*, pages 154–169. Springer, 2014. 1
- [32] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. *SIGGRAPH Comput. Graph.*, 21(4):163–169, aug 1987. 5, 12, 13
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 13
- [34] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 1, 3
- [35] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. 2
- [36] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. *arXiv preprint arXiv:2112.03221*, 2021. 1, 2
- [37] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 3
- [38] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 1, 2, 3, 4
- [39] Norman Müller, Andrea Simonelli, Lorenzo Porzi, Samuel Rota Bulò, Matthias Nießner, and Peter Kontschieder. Autorf: Learning 3d object radiance fields from single view observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3971–3980, June 2022. 1, 3
- [40] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 1, 2, 3
- [41] Charlie Nash, Yaroslav Ganin, SM Ali Eslami, and Peter Battaglia. Polygen: An autoregressive generative model of 3d meshes. In *International Conference on Machine Learning*, pages 7220–7229. PMLR, 2020. 1, 2
- [42] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 2, 3
- [43] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020. 2, 3, 6, 15
- [44] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 2
- [45] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 1, 3
- [46] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. Generating 3d faces using convolutional mesh autoencoders. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 704–720, 2018. 1, 2
- [47] Nikhil Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 1
- [48] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10901–10911, October 2021. 2, 3
- [49] Konstantinos Rematas, Ricardo Martin-Brualla, and Vittorio Ferrari. Sharf: Shape-conditioned radiance fields from a single view. *arXiv preprint arXiv:2102.08860*, 2021. 2, 3
- [50] Katja Schwarz, Axel Sauer, Michael Niemeyer, Yiyi Liao, and Andreas Geiger. Voxgraf: Fast 3d-aware image synthesis with sparse voxel grids. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [51] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-

- aware neural scene representations. *Advances in Neural Information Processing Systems*, 32, 2019. 2, 3, 6, 15
- [52] Jonathan Tremblay, Moustafa Meshry, Alex Evans, Jan Kautz, Alexander Keller, Sameh Khamis, Charles Loop, Nathan Morrical, Koki Nagano, Towaki Takikawa, et al. Rtmv: A ray-traced multi-view synthetic dataset for novel view synthesis. *arXiv preprint arXiv:2205.07058*, 2022. 2, 3
- [53] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European conference on computer vision (ECCV)*, pages 52–67, 2018. 1, 2
- [54] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 3
- [55] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 3
- [56] Xingkui Wei, Zhengqing Chen, Yanwei Fu, Zhaopeng Cui, and Yinda Zhang. Deep hybrid self-prior for full 3d mesh generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5805–5814, 2021. 2
- [57] Mason Woo, Jackie Neider, Tom Davis, and Dave Shreiner. *OpenGL Programming Guide: The Official Guide to Learning OpenGL, Release 1*. Addison-wesley, 1998. 4
- [58] Zhirong Wu, S. Song, A. Khosla, Fisher Yu, Linguang Zhang, Xiaou Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shapes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2015. 2
- [59] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020. 3
- [60] Jun Yang, Yizhou Gao, Dong Li, and Steven L Waslander. Robi: A multi-view dataset for reflective objects in robotic bin-picking. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9788–9795. IEEE, 2021. 3
- [61] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 206–215, 2018. 1, 2
- [62] Shunyu Yao, Tzu Ming Hsu, Jun-Yan Zhu, Jiajun Wu, Antonio Torralba, Bill Freeman, and Josh Tenenbaum. 3d-aware scene manipulation via inverse graphics. In *Advances in neural information processing systems (NIPS)*, pages 1891–1902, 2018. 1
- [63] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1790–1799, 2020. 3
- [64] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2
- [65] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 1, 3
- [66] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 1, 2, 3, 5, 6, 7, 8, 14, 21, 22
- [67] Ye Yu and William AP Smith. Inverserendernet: Learning single image inverse rendering. *arXiv preprint arXiv:1811.12328*, 2018. 1
- [68] Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. Ners: Neural reflectance surfaces for sparse-view 3d reconstruction in the wild. *Advances in Neural Information Processing Systems*, 34, 2021. 1, 2
- [69] 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. 7
- [70] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 3

## A. Detailed Formulations

### A.1. Sparse Convolutions

Sparse convolutions are a variant of standard convolutions that are used in deep learning. In a sparse convolution, only a subset of the input elements is used in the computation, which allows for more efficient use of computation resources and can improve the performance of the convolutional neural network. To perform a sparse convolution, we first define a set of indices that specify which input elements should be used in the convolution. This set of indices is called the "support" of the convolution. We then use these indices to select the relevant input elements and compute the convolution using these elements. This is typically done by applying a filter to the selected input elements and summing the results to produce the output of the convolution.

In the simplest 1 D case, let  $x$  be the input tensor,  $w$  be the convolutional filter, and  $c$  be the support of the convolution (i.e. the set of indices specifying which elements of  $x$  should be used in the convolution). The output of the sparse convolution,  $y$ , can be computed as:  $y = x[c] * w$  where  $*$  denotes the convolution operation, and  $x[c]$  is the subset of elements from  $x$  specified by the support  $c$ . This equation applies the convolutional filter  $w$  to the selected input elements and sums the results to produce the output of the convolution. For more detailed formulation and implementation of the Sparse convolutions we used in our work, please refer to MinkowskiNetwork [9].

## B. Detailed Setup

### B.1. SPARF Dataset

All the rendered images are of  $400 \times 400$  resolution with 4 channels (RGB + alpha channel for background). SPARF has three main splits for every 3D shape: training views (400 views), test views (20 views), and an OOD "hard" views (10 views) as can be shown in Figure 22. Regarding the collected SRFs, Plenoxels [14] is used as the base module. The spherical harmonics dimension of the whole SRFs is  $d = 4 \times 3 = 12$ , while for partial SRFs, it is  $d = 1 \times 3 = 3$ . Most of the hyperparameters used in optimizing the SRFs are the default ones proposed in the Plenoxels paper [14] (as can be seen in the attached code under Svox2/opt/opt-py). However, the following hyperparameters were engineered in order to scale up the optimization and maintain the quality of the SRFs (as can be seen in Figure 23, and 24). Running Plenoxels [14] for fewer iterations ( $3 \times 12K$ ) reduces the time by 30% while maintaining the same PSNR. Using 400 views/shapes in SPARF to optimize the SRFs keep the time manageable in optimization ( $\sim 4$  minutes for the 512 resolution) while maintaining high PSNR ( $\sim 30$ dB). When saving the SRFs, we only save the set of coordinates (integers) and float features (densities and radiance components).

SRF Type	Voxel Resolution	Nb. of Variants	Nb. of SRFs
Partial	32	$4 \times 1$ -view	158,816
		$4 \times 3$ -view	158,816
		$4 \times 1$ -view	158,816
	128	$4 \times 3$ -view	158,816
		$4 \times 1$ -view	158,816
	512	$4 \times 1$ -view	158,816
Whole	32	$1 \times 400$ -view	39,704
	128	$1 \times 400$ -view	39,704
	512	$1 \times 400$ -view	39,704
Total	-	-	1,072,008

Table 6. **SPARF Anatomy.** We show the distribution of the one million SRFs collected in SPARF between multiple resolutions and between whole and partial SRFs.

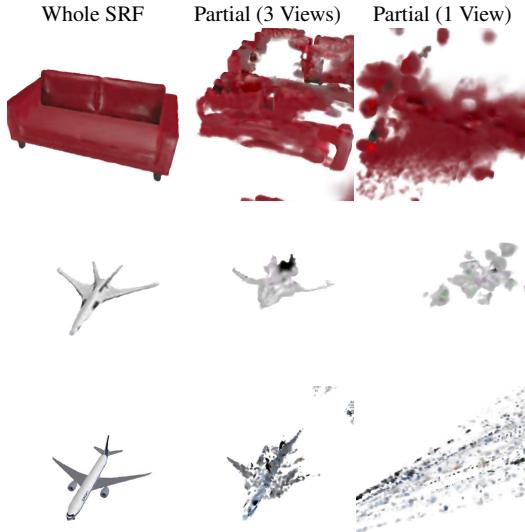


Figure 12. **Whole vs. Partial SRFs.** The partial SRFs are used instead of the few images that generated them as input to the learning pipeline to generate the whole SRFs

The upsampling iteration of Plenoxels is set to  $1 \times 12K$  for faster convergence. The distribution of the collected dataset is detailed in Table 6. More examples of the whole vs. partial SRFs collected in SPARF can be found in Figure 12. The whole SRFs are easily convertible to high-quality meshes using Marching Cubes [32] as shown in Figure 13.

### B.2. SuRFNet Training

We use a voxel resolution of  $128^3$  of the SPARF dataset in most of the learning experiments and visualizations in this work, unless otherwise clearly stated. This choice is to

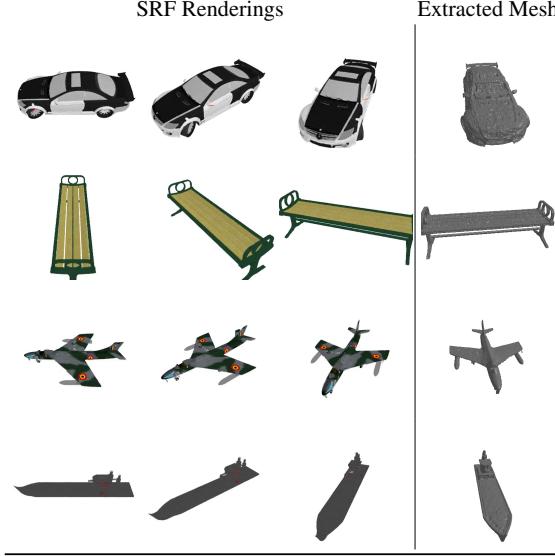


Figure 13. **Extracting 3D Meshes from SRFs.** Since SPARF and SuRFNet live on the 3D voxel’s space, extracting the mesh is straightforward with one pass of MarchingCubes [32].

reduce the computational cost of training the heavy pipeline and to facilitate the development of proper learning methods on SRFs. The input SRF is normalized with a fixed value of 10,0000 for the density and 10 for the colors, to ensure the distribution lies within -1 to 1. The Q-Gaussian std  $\sigma$  is set to 0.444 (studied more in Section C). The strides for the SuRFNet are all set 2, while the network depth  $l = 3$  modules. The batch size for training is 14 when A100 GPUs are used and 6 when V100 GPUs are used. The training saturates at 100 epochs. The optimizer used is AdamW [33] with a learning rate of 0.01, a momentum of 0.9, a weight decay of  $1e - 5$ , and a learning rate exponential decay rate of 0.99. The hyperparameters  $\lambda_R$ ,  $\lambda_\alpha$ ,  $\lambda_\rho$  are all set independently to each class, where a different network is trained on each class separately. Most classes have  $\lambda_\alpha = 30.0$ ,  $\lambda_\rho = 1.0$ ,  $\lambda_R$ . We did not prune the output sparse voxel as this leads to harming performance most of the time and increase the problem of vanishing gradients. The background color of the rendered images  $\mathcal{R}_\phi(\mathbf{F}(\mathcal{X}))$  is masked out from the perceptual loss and the density component  $\alpha$  of the output SRF is also not affected by the perceptual loss, as this can cause excessive densities around the object, leading to deteriorating the SRF output perceptuality. During training with the perceptual loss, three randomly selected images from three different  $\phi$  as used as labels for the three rendered images from the output  $\mathcal{R}_\phi(\mathbf{F}(\mathcal{X}))$ . The SuRFNet is jointly predicting the density and radiance of spherical harmonics colors, but with different heads. More setup details can be found in the attached code and analyzed further in Section C.

For a fair comparison to the baselines PixelNeRF and



Figure 14. **Shiny Objects Corrupts SRFs:** Optimizing SRFs on shiny objects with a reflective material (*left*) results in distorted radiance fields (*right*). These distorted SRFs (of 76 shapes in total) were separated from the main classes in SPARF.

VisionNerf (which use  $64 \times 64$  resolution), we upsample their resolution at inference at test poses while using their pretrained weights of the NMR dataset. The upsampling is using the bicubic sampling of the Pytorch Transforms library. Retraining the methods from scratch on the high resolution  $400 \times 400$  is computationally prohibitive. We train a separate model for each class, to maintain high-quality generation of 3D SRFs.

## C. Additional Analysis

### C.1. Shiny Objects

Some of the rendered objects have reflective materials, resulting in distorted optimized radiance fields for these shapes despite using all of the views. We separate these distorted SRFs (only 76 shapes in total) from the SPARF dataset (see Figure 14).

### C.2. Effect of Dataset Size

We study the effect of increasing the dataset size (Whole SRFs and Partial SRFs) on the generalization performance of SuRFNet in Figure 16,15. It shows that as the dataset size increase (normalized the number of shapes in each class), the generalization performance increase. This scalability effect underlines the importance of SPARF. However, as can be seen from these two figures, partial SRFs scalability is more important than increasing whole SRFs, which justifies collecting 4 variants per resolution (as detailed in Table 6).

### C.3. Loss Ablation Study

For the density threshold  $\alpha_{dense}$  defined in Eq 2 and 3, the validation accuracies of SuRFNet on car class are 13, 14.7, 67.8, 67, 67, 67.2, 62.5 for the values of -0.01, -0.001, 0, 0.001, 0.003, 0.01, 0.03 of  $\alpha_{dense}$  respectively. The hyperparameter  $\sigma$  which governs the spread of the Q-Gaussian loss is studied as follows. the validation accuracies of SuRFNet on airplane class are 52, 70.4, 71.7, 72.1, 72.2, 72.4, 72.3 for the values of 0.1, 0.2, 0.3, 0.4, 0.6, 0.8, and 1.0 of  $\sigma$

Baselines	SPARF Classes													mean
	chair	watercraft	rifle	display	lamp	speaker	cabinet	bench	car	airplane	sofa	table	phone	
Plenoxels [14] (1V)	10.1	12.1	12.6	8.7	14.7	8.7	10.9	11.4	7.7	14.0	9.7	10.5	9.5	10.8
Plenoxels [14] (3V)	10.8	13.3	15.6	9.7	16.2	10.1	12.1	12.1	9.0	15.4	11.4	10.8	10.2	12.1
PixelNeRF [66] (1V)	10.8	14.1	14.2	9.0	15.6	9.2	10.5	12.4	10.1	15.7	11.1	10.6	11.1	11.9
PixelNeRF [66] (3V)	11.0	14.1	14.2	9.3	15.7	9.4	10.6	12.7	10.1	15.7	11.3	10.9	11.4	12.0
VisionNeRF [29] (1V)	16.5	18.4	18.5	15.1	19.3	13.2	16.1	16.3	13.8	21.8	15.1	14.8	14.0	16.4
<b>SuRFNet (ours) (1V)</b>	<b>15.7</b>	<b>15.5</b>	<b>19.1</b>	<b>14.1</b>	<b>18.5</b>	<b>14.5</b>	<b>18.7</b>	<b>15.6</b>	<b>18.1</b>	<b>20.3</b>	<b>16.3</b>	<b>14.1</b>	<b>17.4</b>	<b>16.8</b>
<b>SuRFNet (ours) (3V)</b>	<b>18.6</b>	<b>20.7</b>	<b>20.9</b>	<b>17.1</b>	<b>21.2</b>	<b>18.5</b>	<b>21.7</b>	<b>17.6</b>	<b>18.9</b>	<b>21.9</b>	<b>20.4</b>	<b>16.7</b>	<b>20.0</b>	<b>19.5</b>

Table 7. **SPARF Benchmark on Novel View Synthesis (Normal Test).** We compare the validation PSNR of some of the widely used novel view synthesis techniques on the SPARF dataset for the generalization of novel view synthesis beyond a single example and on the normal testing-views tracks similar to the ones seen in training views. One view (1V) and three views (3V) inputs are reported.

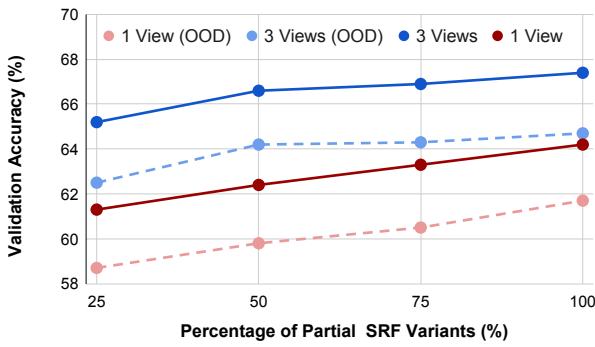


Figure 15. **Scaling-Up Training on SRFs: Partial SRFs.** As the training data (partial SRFs) of radiance fields increase, the generalization improves, as can be seen in the car class here. The 3-view and 1-view metrics are reported with test and OOD metrics.

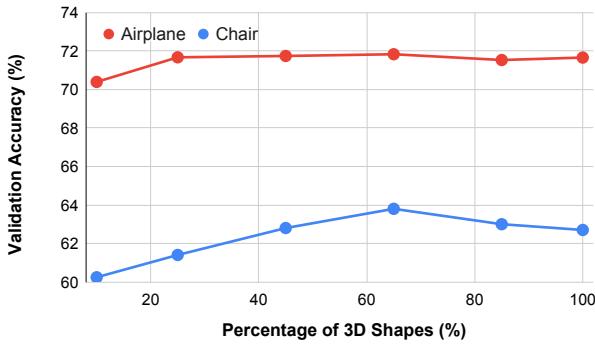


Figure 16. **Scaling-Up Training on SRFs: Whole SRFs.** As the training data of radiance fields increase, the generalization improves across different classes in SPARF.

respectively. The number of coordinates  $\mathbf{c}$  sampled in the Q-Gaussian loss is proportional to the number of coordinates in the input SRFs with multiplier  $K = 40$ . For different values of this multiplier 1, 5, 10, 20, 40, 80, 200, the validation accuracies of SuRFNet trained on airplane class are 72.2,

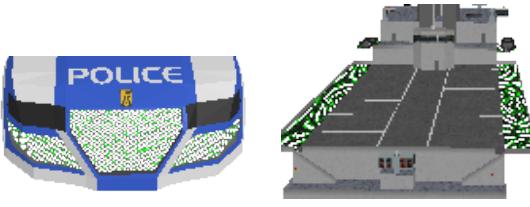


Figure 17. **Rare Cases of Faulty Textures.** Some objects in ShapeNet [6] have doubled textures in some parts, leading to faulty renderings.

output w/o  $L_R$       output w/  $L_R$       whole SRF



Figure 18. **Effect of the Perceptual Loss  $L_R$ .** Adding a perceptual loss on volume-rendered images during training SuRFNet insures the rendered images remain closer to how they should be rendered, as the 3D radiance colors supervision won't guarantee the rendering quality. (*left*): without perceptual loss , (*middle*): with the loss.

72.7, 72.9, 72.5, 71.8, 71.7, and 71.6 respectively.

#### C.4. Faulty Textures

In some rare instance of the shapes in ShapeNet [6], some objects have doubled textures in some areas. This occurs in less than 1% of the data and leads the renderer to render the background instead in these areas (highlighted with green). See Figure 17 for examples of these cases.

#### D. Additional Results

Additional results of normal test tracks benchmark of SPARF are presented in Table 7. Please see figures 21 and 22 for differences between the normal train/test track and the OOD hard track. More comparisons and generations are provided in Figures 26, 27, 25.

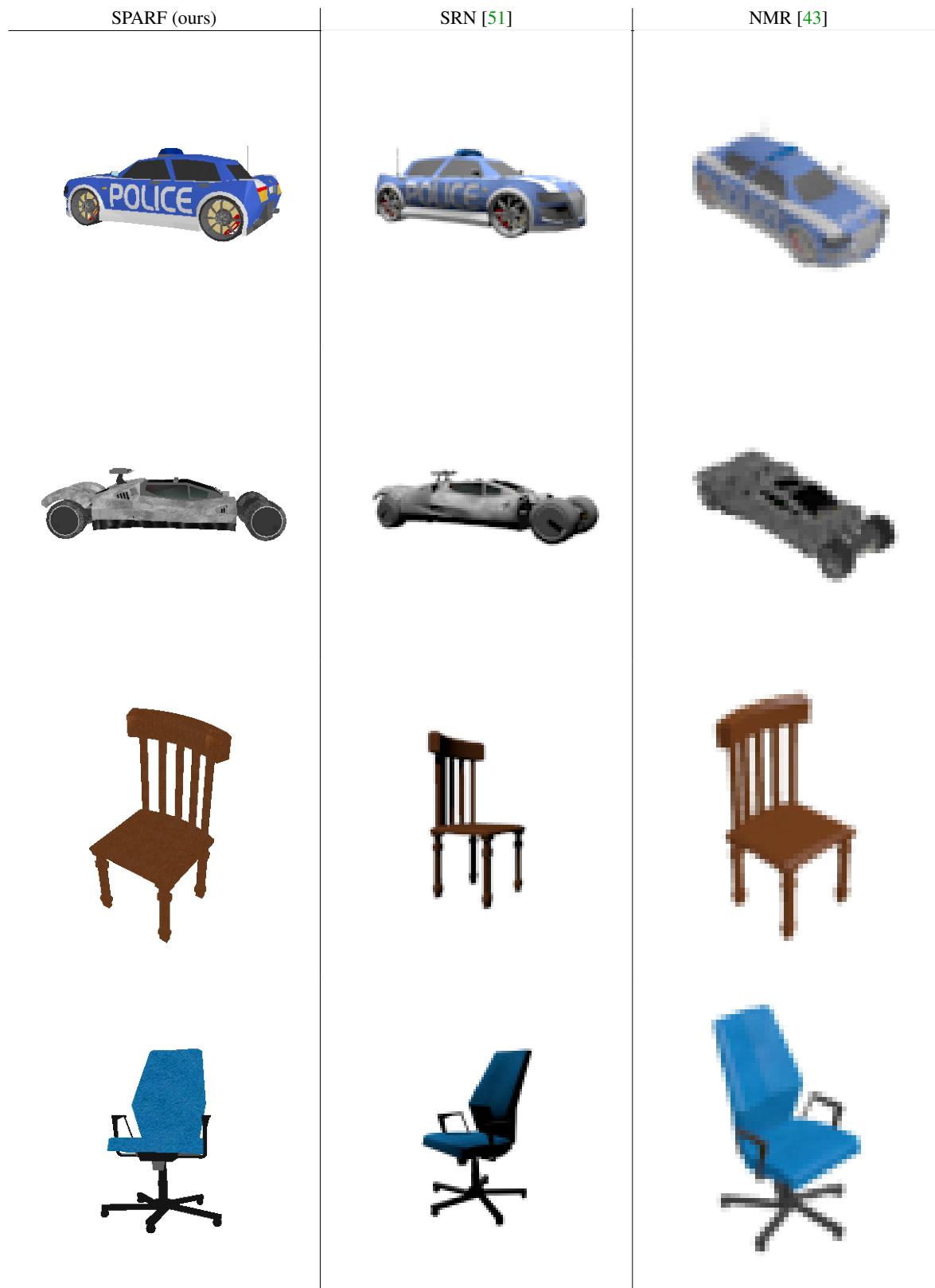


Figure 19. **SPARF vs. other Datasets**. SPARF offers a large-scale high-resolution dataset compared to other posed multi-view datasets. Note that SRN [51] has only cars and chairs, while NMR [43] and SPARF has 13 classes.

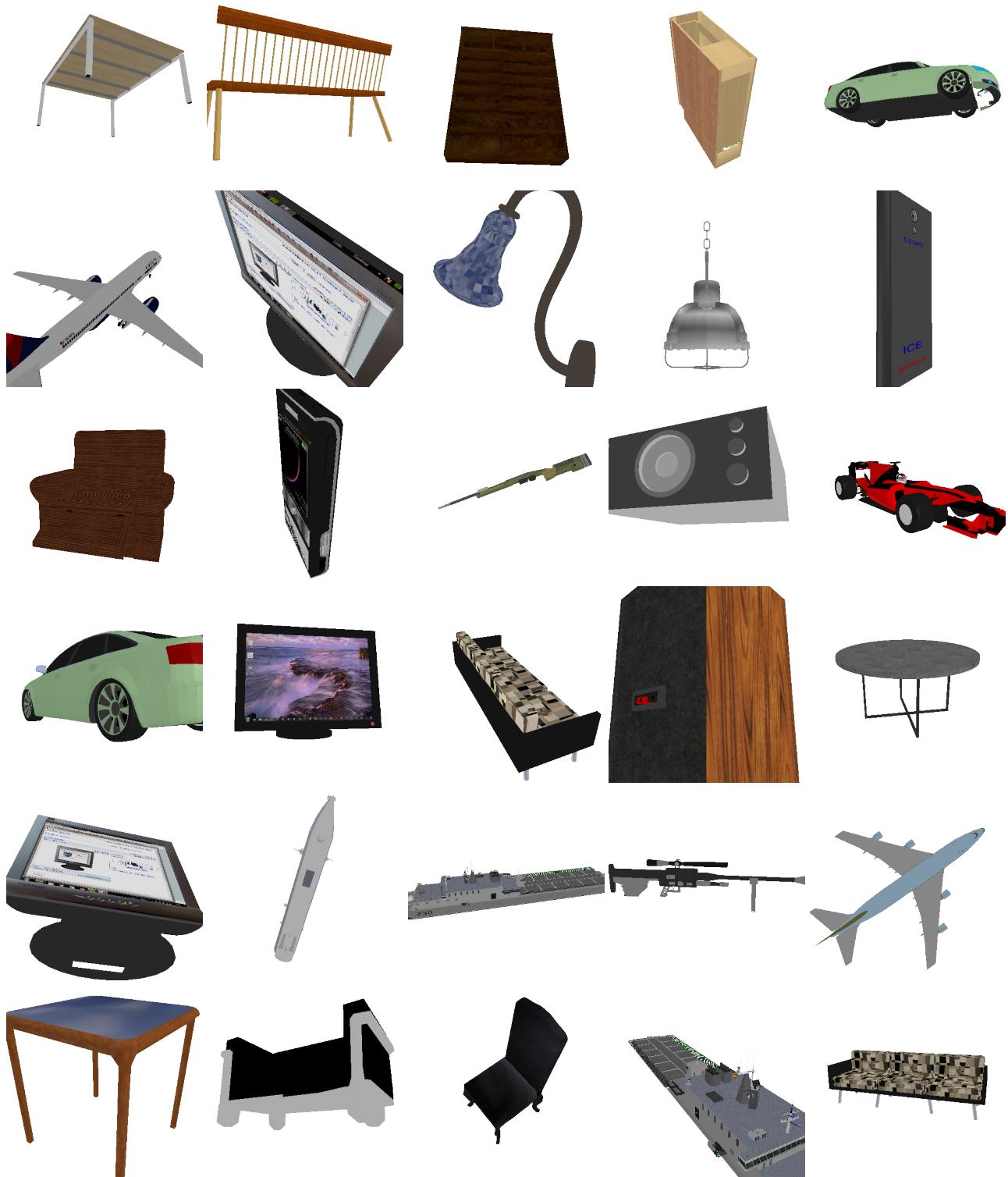


Figure 20. SPARF: a Large Dataset for 3D Shapes Radiance Fields and Novel Views Synthesis.

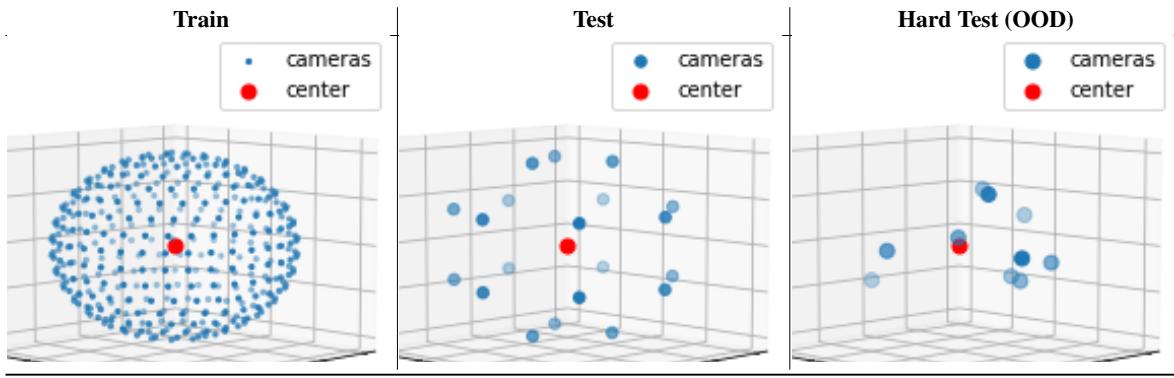


Figure 21. **Cameras Setups for Different SPARF Splits.** Here, we show different visualizations of the camera setups of the three splits of SPARF. (*Train*): 400 deterministic spherical views, (*Test*): 20 random spherical views, (*hard OOD Test*): 10 random views.



Figure 22. **SPARF Splits.** SPARF has three main splits for every 3D shape: training views (400 views), test views (20 views), and OOD “hard” views (10 views) as can be shown in the examples above.

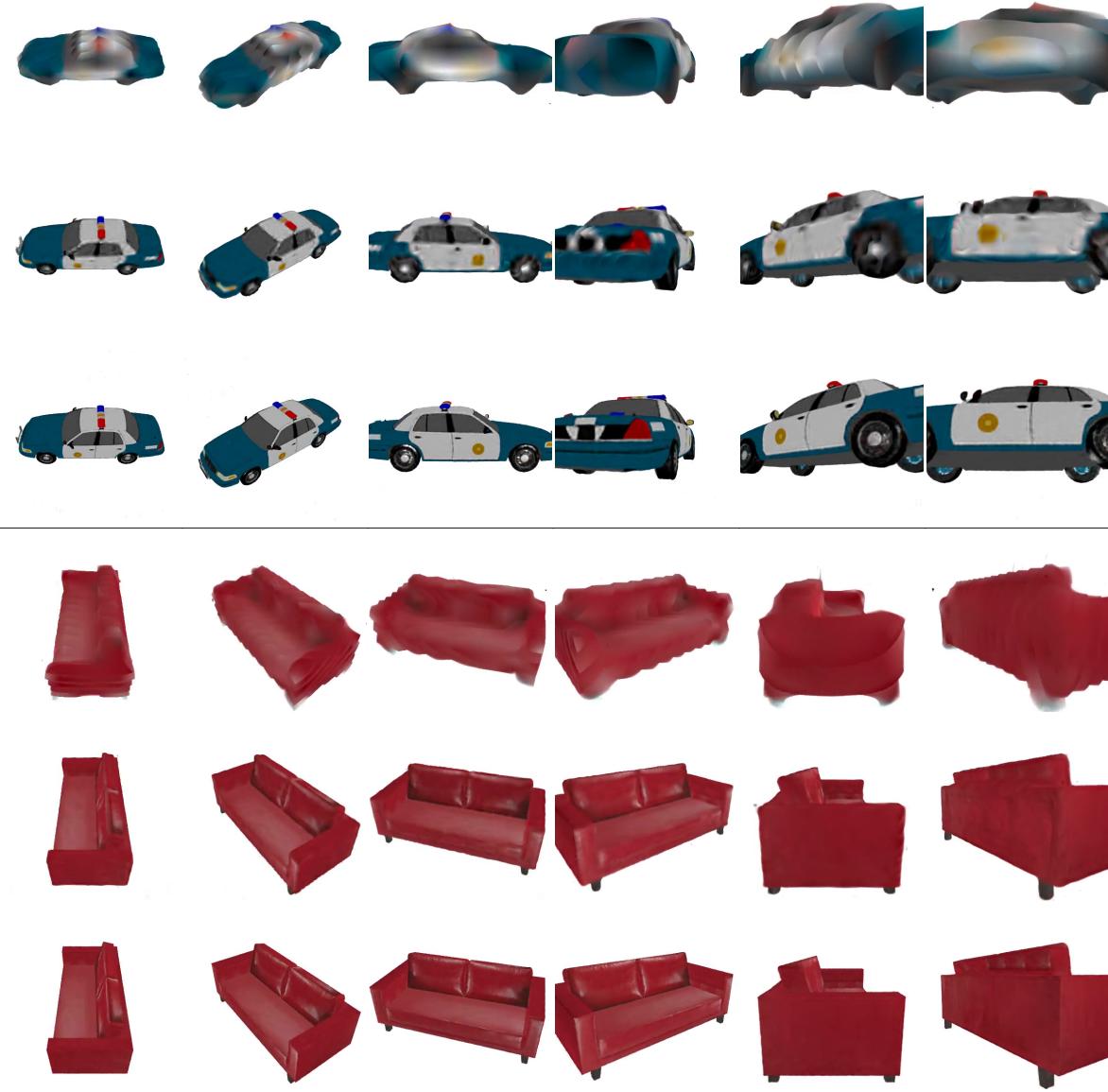


Figure 23. **SRFs: The optimized Sparse Radiance Fields in SPARF 1.** A total of one million SRFs have been collected in SPARF, including on multiple voxel resolutions: 32 (*top*), 128 (*middle*), and 512 (*bottom*) for every 3D shape.

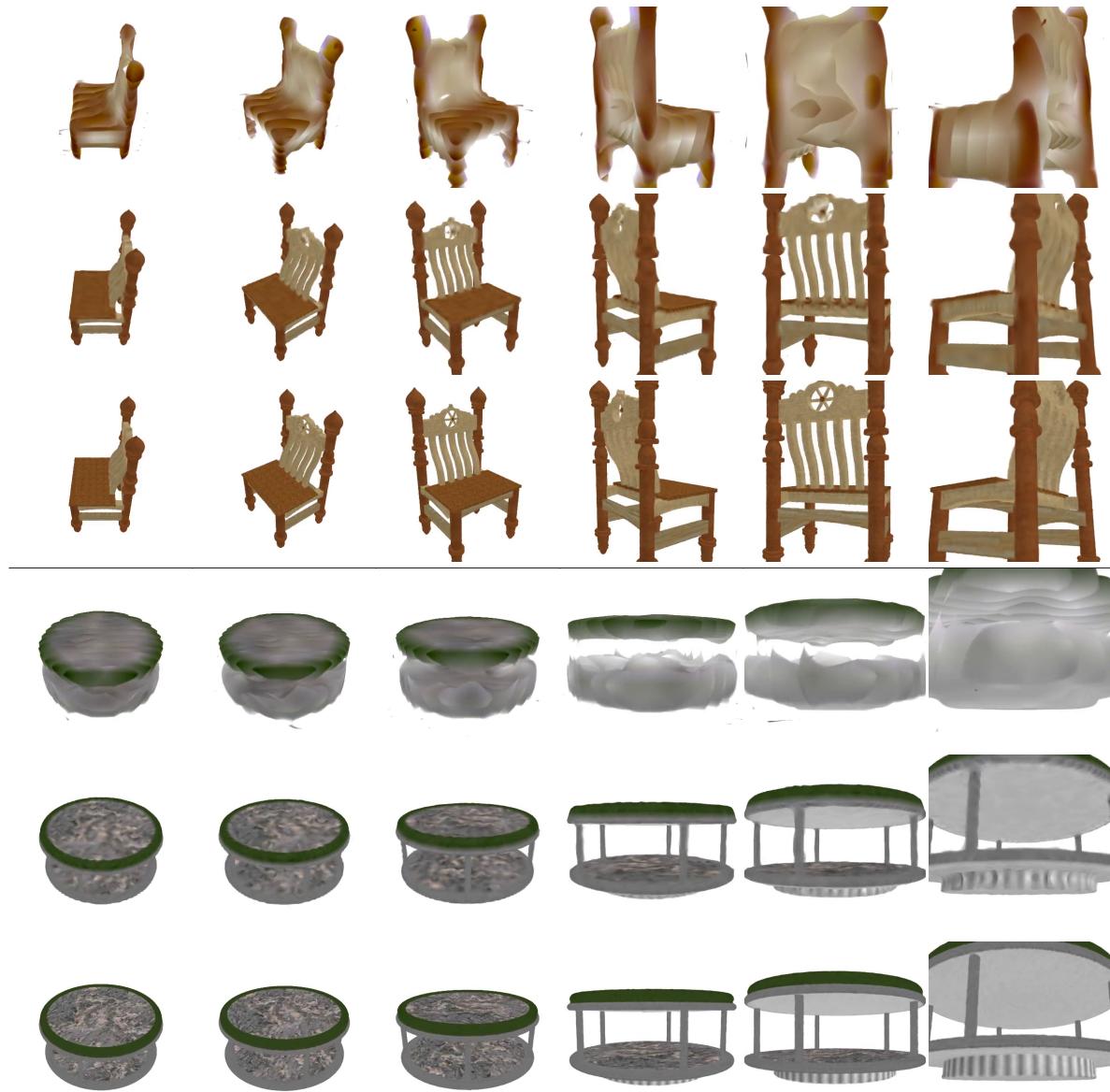
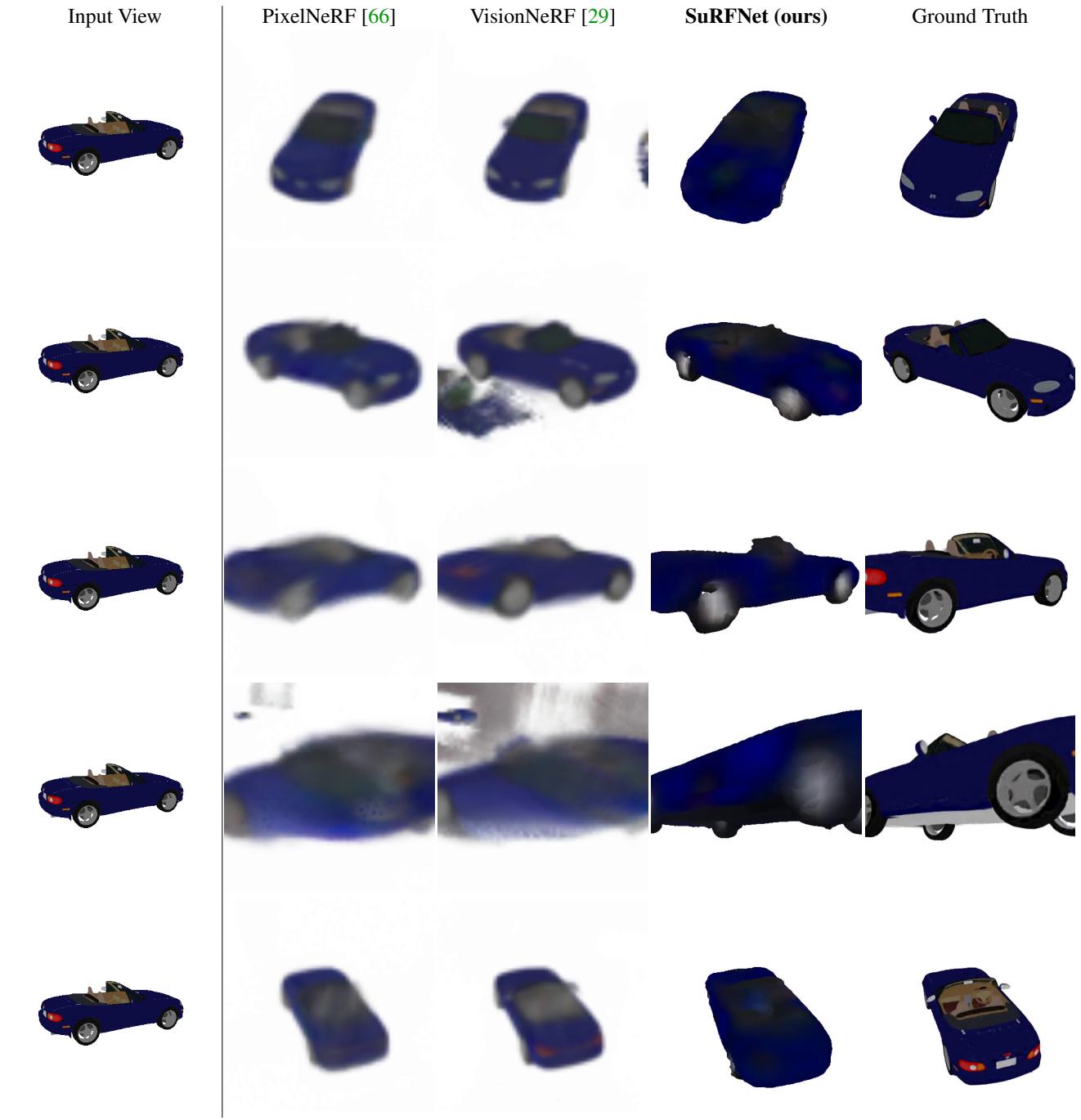


Figure 24. **SRFs: The optimized Sparse Radiance Fields in SPARF 2.** A total of one million SRFs have been collected in SPARF, including on multiple voxel resolutions: 32 (*top*), 128 (*middle*), and 512 (*bottom*) for every 3D shape.



**Figure 25. SuRFNet: Generating High-Resolution Radiance Fields.** We show some volume-rendered sequences based on our SuRFNet voxel radiance field outputs ( $512^3$  resolution), given only 3 images of each shape. This demonstrates the capability of SuRFNet to generate high-resolution sparse voxel SRFs. Note that, here, SURFNet is overfitting on a small dataset in these examples and is not meant for shape generalization.



**Figure 26. Qualitative Comparisons 1.** We show different render from our SuRFNet outputs generated from a single image compared to other methods (pixel-Nerf [66], and VisionNerf [29] ) and whole SRF "GT" renderings. Note that the predicted views are outside the training views distribution (zoomed in randomly). This test highlights the weakness of the 2D-based baselines [29, 66] outside the training track, while our 3D approach maintains multi-view consistency everywhere.



Figure 27. **Qualitative Comparisons 2.** We show different render from our SuRFNet outputs generated from 3 input images compared to other methods (pixel-Nerf [66], and VisionNerf [29] ) and whole SRF "GT" renderings. Note that the predicted views are outside the training views distribution (zoomed in randomly). This test highlights the weakness of the 2D-based baselines [29, 66] outside the training track, while our 3D approach maintains multi-view consistency everywhere.