# SparseCraft: Few-Shot Neural Reconstruction through Stereopsis Guided Geometric Linearization

Mae Younes[†][*][®], Amine Ouasfi[*], and Adnane Boukhayma

Inria, [†]Univ. Rennes, CNRS, IRISA, M2S, France

**Abstract.** We present a novel approach for recovering 3D shape and view dependent appearance from a few colored images, enabling efficient 3D reconstruction and novel view synthesis. Our method learns an implicit neural representation in the form of a Signed Distance Function (SDF) and a radiance field. The model is trained progressively through ray marching enabled volumetric rendering, and regularized with learning-free multi-view stereo (MVS) cues. Key to our contribution is a novel implicit neural shape function learning strategy that encourages our SDF field to be as linear as possible near the level-set, hence robustifying the training against noise emanating from the supervision and regularization signals. Without using any pretrained priors, our method, called SparseCraft, achieves state-of-the-art performances both in novel-view synthesis and reconstruction from sparse views in standard benchmarks, while requiring less than 10 minutes for training. Project page: sparsecraft.github.io

## 1 Introduction

Replicating the 3D world around us digitally in a faithful manner is a long-standing problem that has prompted substantive research in computer vision and graphics alike, with countless downstream applications. While current solutions can provide impressive results, many of them still rely on abundantly informative input, be it in quality (*e.g.* high resolution imagery, depth sensors) or quantity (*e.g.* dense arrays of views). However, due to many constrained scenarios (*e.g.* out-of-the-studio, low budget, *etc.*) and in the interest of wider applicability, the community (*e.g.* [5, 21, 40, 45, 48, 51, 58, 84]) is actively seeking solutions that can deliver under minimal input.

Given a few colored images, we aim to capture both the shape and appearance of the observed object or scene. In practice, we seek metrically accurate 3D reconstruction, and photo-realistic novel view synthesis. In this regard, traditional computational photogrammetry combines structure from motion (SfM) and multi-view stereo (MVS) [61,62] to provide calibration and triangulate an explicit geometry based on matching. However, it can lead to noisy and incomplete meshes in challenging and non Lambertian scenarios. On the other hand, deep learning

---

[*] These authors contributed equally to this work

based implicit neural representations (INR) have emerged as a powerful tool for 3D modelling [43,72]. They have shown ability to learn both detailed shape and radiance though image supervised differentiable volumetric rendering from dense image arrays [72,86]. Learning such accurate implicit shape representations remains challenging when only a few images are available. Current methods in the literature [37,40,58] rely on learned data priors across many training scenes, by conditioning the implicit representation on spatially local features obtained from the sparse input images through generalizable encoders. These can suffer nonetheless from out of distribution generalization issues, and typically require substantial and expensive calibrated multi-view data for training.

Differently, we advocate to fit a neural signed distance (SDF) and radiance functions self-supervisedly to the images. Using a progressively learned hash encoding [44] provides regularization and a more stable and efficient training. We use MVS geometry and color cues to further regularize this challenging learning task. Notice that these cues are readily available, as photogrammetry is typically required to obtain the calibration needed for learning INRs.

Unfortunately, our training is facing noisy labels or phenomena that can be interpreted as such: *e.g.* The MVS geometry can be noisy, and the volumetric rendering supervision can be imprecise due to imperfect calibration, and the inherent bias of geometry based volumetric rendering [10]. To alleviate these challenges, we propose a novel loss rooted in Occam's razor principle. We focus on the surface *i.e.* near the MVS samples, as it is the most critical region in our learning. We hypothesize that excess non-linearity [67] there can lead to overfitting on the noise. Hence, we encourage our SDF to be as linear as possible near the MVS samples, by making the function approximate its first order Taylor expansion, and we integrate the MVS point and normal supervision in this linearization (Section 3.2). We show empirically that this loss leads to considerable improvement in our method, as compared to the previous methods, and also a directly MVS supervised baseline.

We obtain state-of-the-art performances in both novel view synthesis and reconstruction using standard metrics, as-well-as superior qualitative results to previous methods, without using any pre-learned priors, and within shorter training times. In summary, our main contributions are:

• Few-shot 3D reconstruction and novel view synthesis without any pre-learned data priors.

• Leveraging a progressive multi-resolution hash learning strategy in this context.

• A framework, that we call SparseCraft, for harnessing all MVS data: points, normals to regularize the SDF, and color to regularize the diffuse.

• Our novel Taylor expansion inspired losses to regularize SDF learning from sparse multi-view imagery.

## 2   Related Work

**Multi-View Stereo** Conventional MVS can be classified based on scene representation: volumetric [28,29,63], point cloud [11,30], and depth map based [12,62,82].
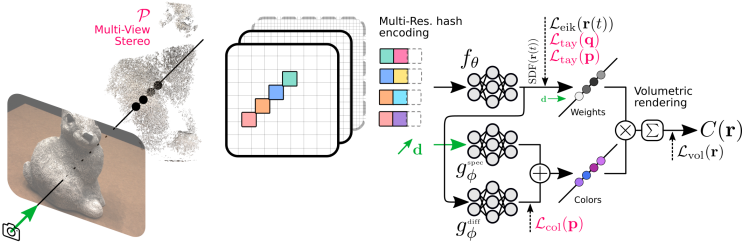
Recently, there has been a preference for depth map based methods due to their versatility, dividing the problem into depth map estimation and fusion stages [12, 62]. Subsequently, Poisson reconstruction [24] is applied to the fused point cloud to produce a watertight mesh. Despite significant advancements, extreme scenarios such as minimal input can still prove challenging for such methods, often resulting in incomplete and inaccurate reconstructions. Nevertheless, we demonstrate that by judiciously leveraging the incomplete fused point cloud generated from one such method (COLMAP [61]), our approach surpasses the state-of-the-art in multi-view few-shot reconstruction.

**Implicit Neural Representations** Implicit Neural fields employ deep neural networks to model 2D or 3D data as continuous functions, overcoming many of the limitations of explicit ones (*e.g.* meshes [19,23,71] and point clouds [1,9,25]) in modelling shape, radiance and light fields (*e.g.* [3,16,31,32,43,72,86]). The seminal work (NeRF) [43], which combines volume rendering and implicit representations, has paved the way for learning diverse tasks, including novel view synthesis [41,43], 3D generation [17,54], deformation [52,55,57], and video rendering [8,34,36,80]. More recently, attention was shed on implicit surface reconstruction, with a focus on single-stage optimization and robust representation potential [83,87]. This was improved through novel weight functions involving SDFs for color accumulation during volumetric rendering [72,86]. However, persistent challenges such as geometric bias arising from discrete sampling and other factors [10,92] still remain.

Efforts have been directed towards addressing the time-intensive training associated with these methods. *i.e.*, NeRF-based work [44,60,68] introduced voxel-grid features. Subsequently, other literature [74,78] extended them to surfaces. Lately, Neuralangelo [35] proposed leveraging multi-resolution hash grids with numerical gradient computation and a topology warm-up strategy for neural surface reconstruction. While achieving high-fidelity geometry from dense images, it comes with a considerable training time cost. Inspired by the latter, our work leverages numerical gradients and hash encoding, and additionally employs an occupancy grid for sampling, to strike a balance between reconstruction quality and training speed.

**Novel-View from Sparse Input** Existing work has tackled this task by incorporating additional information, such as normalization-flow [45], perceptual [90] and diffusion-based [79] regularization, depth supervision [7,59,76], and enforcing cross-view semantic consistency [18]. Conversely, another line of work [4,6,88] strives to develop transferable models by training on a large, curated dataset, eschewing the use of external models. Recent investigations posit that geometry is a pivotal factor in few-shot neural rendering, advocating for geometry regularization [45] to enhance performance. However, these methods require resource-intensive pre-training on tailored multi-view datasets [4,6,88] or employing costly training-time patch rendering [18,45], thus introducing significant overhead.

Conversely, other approaches propose regularization strategies during single scene fitting. These include frequency encoding regularization [84], entropy con-

**Fig. 1:** Overview: In this toy example, we illustrate inference given 4 samples $\{\mathbf{r}(t)\}$ on a ray $\mathbf{r}$ (where the last hash resolution is not active yet). Dashed arrows symbolize losses operating mid-training. SparseCraft leverages differentiable volumetric rendering to learn a SDF based implicit representation given a few images, using MVS cues as regularization ( losses in Red).

straints on density [26], utilizing a mixture density model [65] or exploiting flipped reflection rays as augmentation [64]. In this work, we demonstrate that explicit regularization, aimed at linearizing the signed distance function in proximity to the surface with guidance from an incomplete point cloud derived from a classical MVS method, along with the incorporation of a progressive hash encoding, not only enhances the surface reconstruction of our SDF based method, but also enables it to surpass state-of-the-art NeRF-based approaches in rendering quality in the few-shot object-centric setting.

**Reconstruction from Sparse Input** For this task, geometric priors [10, 50, 70, 89, 91] have been proposed to enhance reconstruction in the single scene fitting setting. However, these methods are slow to train and still display artifacts and failures. Generalizable novel view synthesis models [4, 39, 53, 66, 69, 73, 88] can be repurposed for reconstruction by carefully adjusting the density threshold for extraction. However, their reconstructions tend to be noisy and not as robust as reconstruction methods. Generalizable surface reconstruction methods (from images [40, 58] as well as point clouds [2, 47, 49, 53]) are still prone to failure for out-of-distribution scenes/views. Another noteworthy work [77] employs the neural rendering of an implicit reconstruction method to improve the MVS performance of deep MVS models in the few-shot setting. In contrast to all the aforementioned work, our rapidly trained method achieves state-of-the-art results for the few-shot reconstruction task without relying on pre-learned priors.

## 3    Method

Given a few input colored images $\{I_i\}_{i=1}^N$, our goal is to recover the shape and appearance of the observation. We achieve this by learning implicit shape and radiance functions simultaneously. We model the shape $\mathcal{S}$ with a SDF $f$ parameterized with a neural network $f_\theta$. We model the radiance as a view direction $\mathbf{d}$ and location $\mathbf{x}$ dependent 3-channel color function $g$ parameterized through a neural network $g_\phi$, *i.e.* $g(\mathbf{x}, \mathbf{d}) = \mathbf{c}$ where $\mathbf{c} \in [0, 1]^3$. The inferred

shape $\hat{S}$ can be obtained at test time as the zero level set of the learned SDF $f_\theta$ at convergence: $\hat{S} = \{\mathbf{x} \in \mathbb{R}^3 \mid f_\theta(\mathbf{x}) = 0\}$. Concurrently, given a target new view point, a novel image $\hat{I}$ can be generated through ray-wise volumetric rendering [22] per pixel, using the converged neural SDF and radiance fields $f_\theta$ and $g_\phi$ respectively.

### 3.1   Learning Implicit Neural Shape and Radiance

The NeRF [43] framework enables learning a volumetric scene representation through a synthesis and compare procedure between generated and ground-truth pixel values. Let us assume a ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$, where $\mathbf{o}$ is the camera origin and $\mathbf{d}$ the ray direction. The color $C$ of the pixel corresponding to a ray $\mathbf{r}$ can be generated through integration along the ray:

$$C(\mathbf{r}) = \int T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt. \tag{1}$$

where $\sigma(\mathbf{r}(t))$ denotes volume density, which represents a differential opacity signaling the amount of radiance accumulated by a ray passing through the point $\mathbf{r}(t)$. $T(t)$ denotes transparency, *i.e.* the accumulated transmittance along the ray until $t$, which can be derived from density accordingly:

$$T(t) = \exp\left(-\int^t \sigma(\mathbf{r}(s))ds\right). \tag{2}$$

As extracting geometry by thresholding density $\sigma$ yields suboptimal and noisy results, recent literature proposed to involve a SDF $f$ in the volumetric rendering equation, by defining a function $\Psi$ that transforms the SDF into density $\sigma(\mathbf{r}(t))$ for Yariv *et al.* [86], the weighting function $T(t)\sigma(\mathbf{r}(t))$ in the case of Wang *et al.* [72], and most recently the transparency $T(t)$ in the work by Wang *et al.* [75]. We follow here the latter representation, where the transformation $\Psi$ is chosen to satisfy monotony and boundary conditions fit for $T(t)$:

$$T(t) = \Psi_s(f(\mathbf{r}(t))) = \frac{1}{1 + \exp(-sf(\mathbf{r}(t)))}, \tag{3}$$

where $s$ controls the slope of the transformation. In practice, the integral in Equation 1 is approximated using discrete samples $\{t_i\}$ with the quadrature rule [42]. Using our SDF and radiance neural networks $f_\theta$ and $g_\phi$, the inferred color of a ray then writes:

$$C(\mathbf{r}) = \sum T_i\left(1 - e^{-\sigma_i(t_{i+1}-t_i)}\right)g_\phi(\mathbf{r}(t_i), \mathbf{d}), \tag{4}$$

where the transparency and density are obtained from the SDF network [75]. Model parameters $\theta$ and $\phi$ can be optimized at this stage using the following empirical risk minimization:

$$\min_{\theta, \phi} \mathop{\mathbb{E}}_{\substack{\mathbf{r} \sim \mathcal{R} \\ t \sim \mathcal{T}_\mathbf{r}}} \mathcal{L}_{\mathrm{vol}}(\mathbf{r}) + \mathcal{L}_{\mathrm{eik}}(\mathbf{r}(t)), \tag{5}$$

where $\mathcal{R}$ symbolizes a distribution over training rays among all training images, and $\mathcal{T}_{\mathbf{r}}$ is the sample distribution along a ray $\mathbf{r}$. $\mathcal{L}_{\text{vol}}$ is the photometric reconstruction loss based on discretized volumetric rendering (Equation 4), while $\mathcal{L}_{\text{eik}}$ is the Eikonal regularization [13] that helps avoid the all zero SDF degenerate solution:

$$\mathcal{L}_{\text{vol}}(\mathbf{r}) = ||C(\mathbf{r}) - C_{gt}(\mathbf{r})||_1. \tag{6}$$

$$\mathcal{L}_{\text{eik}}(\mathbf{r}(t)) = \left(||\nabla f(\mathbf{r}(t))||_2 - 1\right)^2. \tag{7}$$

### 3.2   Regularization with Stereopsis Cues

Learning a SDF and radiance field conjointly from few images can be an underconstrained problem, which underpins the need for additional regularization. Fitting 3D INRs to images entails typically an automatic calibration preprocess. The latter estimates camera intrinsic and extrinsic parameters, which are key to performing 3D consistent volumetric rendering and/or ray marching. The main method of choice in this context remains COLMAP [61,62] (SfM + MVS). Hence, without additional overhead, we can acquire the dense fused MVS point cloud, with its point-wise normal estimations and color labels using the sparse input images. We propose subsequently to regularize our training with these additional cues. We note that while MVS points have been exploited before in learning NeRFs from sparse [7,59,76] and dense images [81], and SDF based radiance from dense images [21], we propose differently here to use these cues in learning SDF based radiance in the sparse setting. Additionally, and to the best of our knowledge, our work is also the first to suggest leveraging the color and normal MVS labels, and not only the point spatial locations.

However, the MVS surface samples come with a considerable deal of noise, while also being incomplete, due to inaccuracies in the matching and triangulation process that further intensify in our sparse input setting, along with MVS related limitations when dealing with challenging surfaces (*e.g.* textureless and reflective surfaces). Furthermore, we can argue that even our volumetric rendering based supervision is prone to noise. As a matter of fact, this noise can be manifested in *e.g.* 3D inconsistent supervision emanating from imprecision in the calibration, and also in the inherent bias [81] arising from volumetric integration based geometry modeling, as opposed to *e.g.* a pinpoint root rasterization based geometry modelling. We propose a novel strategy to remedy these challenges in the following.

**Taylor Expansion Based Geometric Regularization** We focus on the level set of our SDF, where the most crucial knowledge for rendering concentrates. We hypothesize that encouraging our SDF to be as linear as possible there can robustify it against the noise introduced above as intuitively, overly complex models are more likely to overfit on noisy samples [56]. We derive a loss that can achieve this linearization efficiently, while integrating MVS point and normal label supervision seamlessly.

We denote by $\mathcal{P} \subset \mathbb{R}^{3 \times M}$ the MVS point cloud obtained from input images $\{I_i\}$. We note that each sample $\mathbf{p} \in \mathcal{P}$ comes with a normal $\mathbf{n}_{\text{MVS}}(\mathbf{p})$ and color

$\mathbf{c}_{\mathrm{MVS}}(\mathbf{p})$ estimation. We generate a pool of query points near the surface by sampling around the MVS points following a normal distribution, *i.e.* $\{\mathbf{q} \sim \mathcal{N}(\mathbf{p}, \sigma_\epsilon \mathbf{I}_3)\}$ where the standard deviation $\sigma_\epsilon$ decreases proportionately with the progressive learning step $\epsilon$ during training. We recompute subsequently the nearest point $\mathbf{p}$ in $P$ for each sample $\mathbf{q}$, thus forming the following set of training pairs:

$$\mathcal{Q} := \{(\mathbf{q}, \mathbf{p}), \mathbf{p} = \min_{\mathbf{v} \in \mathcal{P}} ||\mathbf{v} - \mathbf{q}||_2\}. \tag{8}$$

Given a pair $(\mathbf{q}, \mathbf{p})$ in $\mathcal{Q}$, let us consider the first order Taylor polynomial approximation of our SDF $f_\theta$ around $\mathbf{q}$, and evaluate it at $\mathbf{p}$, secure in the knowledge that it is in the direct vicinity of $\mathbf{q}$:

$$f_\theta(\mathbf{p}) \approx f_\theta(\mathbf{q}) + \nabla f_\theta(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}). \tag{9}$$

Leveraging the constraint that points $\mathbf{p}$ need to belong to the zero level set, we can derive the following loss:

$$\mathcal{L}_{\mathrm{tay}}(\mathbf{q}) = ||f_\theta(\mathbf{q}) + \nabla f_\theta(\mathbf{q})^\top (\mathbf{p} - \mathbf{q})||_2. \tag{10}$$

Note that this loss encourages both our function to have minimal curvature near the surface, and MVS points to coincide with the level set.

Multiplying by gradient $\nabla f_\theta$ and rearranging Equation 9 leads to the following approximation:

$$\mathbf{p} \approx \mathbf{q} - f_\theta(\mathbf{q}) \cdot \frac{\nabla f_\theta(\mathbf{q})}{||\nabla f_\theta(\mathbf{q})||_2^2}. \tag{11}$$

Hence, our loss $\mathcal{L}_{\mathrm{tay}}(\mathbf{q})$ can also be interpreted as supervising a single step of Newton root finding on the SDF $f_\theta$, initialized at $\mathbf{q}$, with its nearest neighbor in the MVS point cloud $\mathbf{p}$.

Let us consider now the Taylor approximation of our SDF around $\mathbf{p}$ conversely, as evaluated at query $\mathbf{q}$:

$$f_\theta(\mathbf{q}) \approx f_\theta(\mathbf{p}) + \nabla f_\theta(\mathbf{p})^\top (\mathbf{q} - \mathbf{p}). \tag{12}$$

We can leverage here the additional constraint that the normalized gradient of the SDF needs to approximate the surface normal, *i.e.* $\mathbf{n}_{\mathrm{MVS}}(\mathbf{p}) \approx \nabla f_\theta(\mathbf{p})/||\nabla f_\theta(\mathbf{p})||_2$. Thus, we can derive the loss:

$$\mathcal{L}_{\mathrm{tay}}(\mathbf{p}) = ||f_\theta(\mathbf{q}) - ||\nabla f_\theta(\mathbf{p})||_2 \cdot \mathbf{n}_{\mathrm{MVS}}(\mathbf{p})^\top (\mathbf{q} - \mathbf{p})||_2. \tag{13}$$

Table 3 and Figure 7 show the benefit of using these Taylor losses as opposed to standard direct supervision.

**Color Regularization** The color labels provided by MVS are averaged from the input images, so we propose to use them as a supervision to the diffuse component of our radiance. Hence, as illustrated in Figure 1, and following [35], our color network consists of two small MLPs $g_\phi^{\mathrm{diff}}$ and $g_\phi^{\mathrm{spec}}$ modelling view independent and view dependent radiance respectively:

$$g_\phi(\mathbf{x}, \mathbf{d}) := g_\phi^{\mathrm{spec}}(\mathbf{x}, \mathbf{F}_\theta(\mathbf{x}), \nabla f_\theta(\mathbf{x})/||\nabla f_\theta(\mathbf{x})||_2, \mathbf{d}) + g_\phi^{\mathrm{diff}}(\mathbf{x}, \mathbf{F}_\theta(\mathbf{x})), \tag{14}$$

where $\mathbf{F}_\theta$ is a feature extracted from the geometry network $f_\theta$. Our color regularization applied at MVS points then writes:

$$\mathcal{L}_{\text{col}}(\mathbf{p}) = ||g_\phi^{\text{diff}}(\mathbf{p}) - \mathbf{c}_{\text{MVS}}(\mathbf{p})||_1. \tag{15}$$

Finally, we can learn our implicit neural representation through the following combined optimization:

$$\min_{\theta,\phi} \ \mathbb{E}_{\substack{\mathbf{r}\sim\mathcal{R} \\ t\sim\mathcal{T}_\mathbf{r} \\ (\mathbf{q},\mathbf{p})\sim\mathcal{Q}}} \ \mathcal{L}_{\text{vol}}(\mathbf{r}) + \mathcal{L}_{\text{eik}}(\mathbf{r}(t)) + \mathcal{L}_{\text{tay}}(\mathbf{q}) + \mathcal{L}_{\text{tay}}(\mathbf{p}) + \mathcal{L}_{\text{col}}(\mathbf{p}). \tag{16}$$

Figure 1 provides a visual summary of our method.
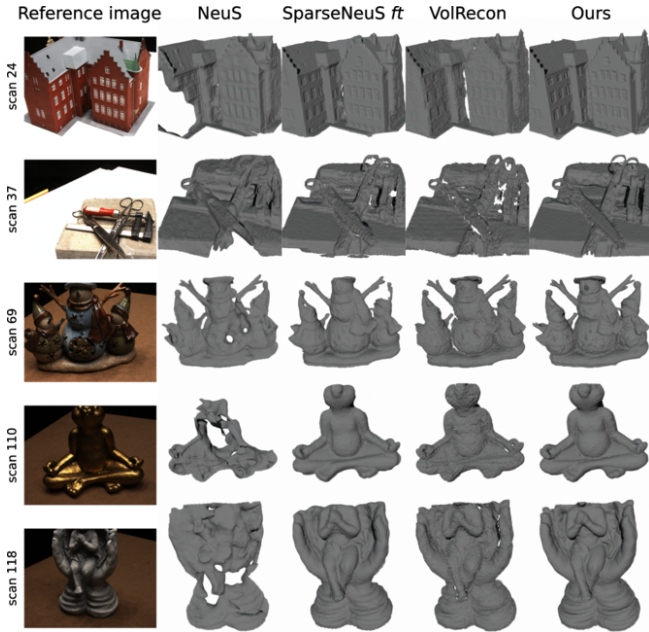
### 3.3   Fast Progressive Learning

Our implementation of the INR network builds on the seminal work in [44]. Our SDF $f_\theta$ consists of an efficiently CUDA implemented multi-resolution hash encoding followed by a small MLP, and the radiance network $g_\phi$ consists of two small MLPs. We use an explicit occupancy grid that guides the sampling along rays (*i.e.* $t \sim \mathcal{T}_\mathbf{r}$) for inference. This combination allows for fast training.

While progressive learning through positional encoding or learnable features was introduced previously for learning NeRFs from dense (*e.g.* [38]) and sparse (*e.g.* [84]) images, and SDF based radiance from dense images (*e.g.* [35]), we propose here to explore this strategy for SDF based radiance learning in the sparse setting for the first time to the best of our knowledge. Differently from Neuralangelo [35], we use the progressive hash encoding to regularize the training in the few shot setting. Hence, it is applied throughout the whole training, rather than its use as a warm-up strategy in [35]. We note that progressively activating hash resolutions during training reduces overfitting and improves the stability of the training in our experiments, and also improves the rendering quality. We also use numerical gradient to approximate derivatives, which allows to back-propagate gradients to more hash cells in the training. The step size of the derivative computation $\epsilon$ is scheduled progressively in concordance with the hash dimensions, as recommended in [35]. Details of the scheduling of our progressive learning are reported in the supplementary material.
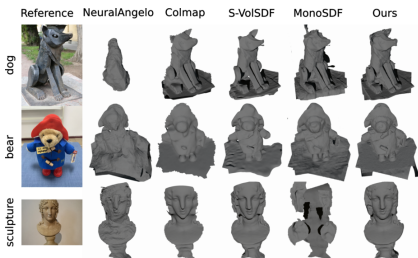
## 4   Implementation Details

We build upon the instant-nsr-pl [15] implementation of Neuralangelo and utilize Nerfacc's [33] accelerated sampling with occupancy grid. Our hash resolution spans from $2^2$ to $2^{11}$ with 32 levels, and we employ a multi-level optimization strategy. We use AdamW optimizer with a learning rate schedule and a combination of losses with varying weights. For more details on the implementation, including the architecture of our MLPs, training protocol, and cues sampling, please refer to the supplementary material.

**Fig. 2:** Qualitative comparison of surface reconstruction in DTU from 3 views. **SparseNeuS and VolRecon use deep data-driven priors, whereas we do not**.
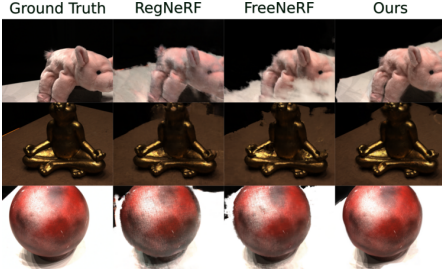


**Fig. 3:** Qualitative comparison of surface reconstruction in BMVS from 3 views.
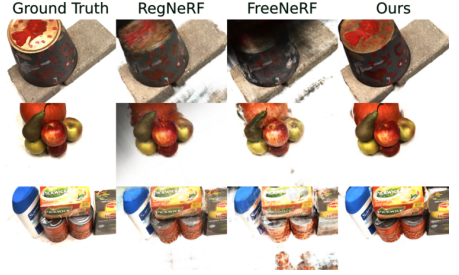


**Fig. 4:** Qualitative comparison of surface reconstruction on T&T from 24 uniformly sampled views.

## 5    Experiments

**Fig. 5:** Qualitative comparison of novel view synthesis in DTU from 3 views.



**Fig. 6:** Qualitative comparison of novel view synthesis in DTU from 6 views.

| Scan | 24 | 37 | 40 | 55 | 63 | 65 | 69 | 83 | 97 | 105 | 106 | 110 | 114 | 118 | 122 | Mean ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| COLMAP [61] | **0.9** | 2.89 | 1.63 | 1.08 | 2.18 | 1.94 | 1.61 | *1.3* | 2.34 | 1.28 | 1.1 | 1.42 | 0.76 | 1.17 | *1.14* | 1.52 |
| IDR [87] | 4.01 | 6.4 | 3.52 | 1.91 | 3.96 | 2.36 | 4.85 | 1.62 | 6.37 | 5.97 | 1.23 | 4.73 | 0.91 | 1.72 | 1.26 | 3.39 |
| VolSDF [86] | 4.03 | 4.21 | 6.12 | *0.91* | 8.24 | *1.73* | 2.74 | 1.82 | 5.14 | 3.09 | 2.08 | 4.81 | 0.6 | 3.51 | 2.18 | 3.41 |
| UNISURF [46] | 5.08 | 7.18 | 3.96 | 5.3 | 4.61 | 2.24 | 3.94 | 3.14 | 5.63 | 3.4 | 5.09 | 6.38 | 2.98 | 4.05 | 2.81 | 4.39 |
| NeuS [72] | 4.57 | 4.49 | 3.97 | 4.32 | 4.63 | 1.95 | 4.68 | 3.83 | 4.15 | 2.5 | 1.52 | 6.47 | 1.26 | 5.57 | 6.11 | 4.00 |
| SparseNeuS ft [40] | 1.29 | 2.27 | *1.57* | 0.88 | 1.61 | 1.86 | *1.06* | 1.27 | 1.42 | 1.07 | 0.99 | 0.87 | 0.54 | *1.15* | 1.18 | *1.27* |
| MVSNeRF [4] | 1.96 | 3.27 | 2.54 | 1.93 | 2.57 | 2.71 | 1.82 | 1.72 | 2.29 | 1.75 | 1.72 | 1.47 | 1.29 | 2.09 | 2.26 | 2.09 |
| PixelNerf [88] | 5.13 | 8.07 | 5.85 | 4.4 | 7.11 | 4.64 | 5.68 | 6.76 | 9.05 | 6.11 | 3.95 | 5.92 | 6.26 | 6.89 | 6.93 | 6.28 |
| SparseNeuS infer [40] | 1.68 | 3.06 | 2.25 | 1.1 | 2.37 | 2.18 | 1.28 | 1.47 | 1.8 | 1.23 | 1.19 | 1.17 | 0.75 | 1.56 | 1.55 | 1.64 |
| VolRecon [58] | 1.2 | 2.59 | 1.56 | 1.08 | *1.43* | 1.92 | 1.11 | 1.48 | 1.42 | 1.05 | 1.19 | 1.38 | 0.74 | 1.23 | 1.27 | 1.38 |
| ReTR [37] | 1.05 | *2.31* | **1.44** | 0.98 | **1.18** | **1.52** | 0.88 | 1.35 | **1.3** | 0.87 | *1.07* | 0.77 | *0.59* | 1.05 | 1.12 | 1.17 |
| Ours (SparseCraft) | *1.17* | **1.74** | 1.8 | **0.7** | 1.19 | 1.53 | 0.83 | 1.05 | 1.42 | 0.78 | 0.8 | 0.56 | 0.44 | 0.77 | 0.84 | 1.04 |

**Table 1:** Quantitative results of sparse view surface reconstruction on 15 testing scenes of DTU dataset [20]. We report Chamfer distance (lower is better). Best scores are in **bold**, second best are in underlined and third best are in *italic*.

## 5.1 Datasets and Setups

We follow the experimental settings in RegNeRF [45] for novel view synthesis from sparse views (3, 6 and 9). We follow SparseNeuS [40] for reconstruction from 3 views. More details can be found in the supplementary material. For all qualitative figures, more examples are provided in the supplementary material.

**Few-Shot Novel View Synthesis** We evaluate on 15 scenes from the DTU dataset [20]. We follow the protocol in sparse NeRF-based methods [45,64,84,88]. Following these methods, we report here the foreground metric, and provide the full image metric in the supplementary material as-well. We report PSNR, SSIM, VGG LPIPS scores, and the geometric average, following [45].

We relay results reported by [65,79,84]. We mainly compare against state-of-the-art generalizable methods PixelNeRF [88], Stereo Radiance Fields (SRF) [6] and MVSNeRF [4] as pretrained and fine-tuned (denoted with "ft") by RegN-eRF [45]. We also compare against NeRF-based methods that use external priors [18,45,79], as well as NeRF-based regularization methods [64,65,84].

**Few-Shot Reconstruction** We evaluate on datasets DTU [20], BlendedMVS [85] and Tanks & Temples [27]. We use the same 15 testing scenes as SparseNeuS [40] (Please note that the DTU splits and views used for novel view synthesis and

| | Object PSNR ↑ | | | Object SSIM ↑ | | | Object LPIPS ↓ | | | Object Average ↓ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 3 views | 6 views | 9 views | 3 views | 6 views | 9 views | 3 views | 6 views | 9 views | 3 views | 6 views | 9 views |
| SRF [6] | 15.32 | 17.54 | 18.35 | 0.671 | 0.73 | 0.752 | 0.304 | 0.25 | 0.232 | 0.171 | 0.132 | 0.12 |
| PixelNeRF [88] | 16.82 | 19.11 | 20.4 | 0.695 | 0.745 | 0.768 | 0.27 | 0.232 | 0.22 | 0.147 | 0.115 | 0.1 |
| MVSNerf [4] | 18.63 | 20.7 | 22.4 | *0.769* | 0.823 | 0.853 | 0.197 | 0.156 | 0.135 | 0.113 | 0.088 | 0.068 |
| SRF ft [6] | 15.68 | 18.87 | 20.75 | 0.698 | 0.757 | 0.785 | 0.281 | 0.225 | 0.205 | 0.162 | 0.114 | 0.093 |
| PixelNeRF ft [88] | 18.95 | 20.56 | 21.83 | 0.71 | 0.753 | 0.781 | 0.269 | 0.223 | 0.203 | 0.125 | 0.104 | 0.09 |
| MVSNeRF ft [4] | 18.54 | 20.49 | 22.22 | *0.769* | 0.822 | 0.853 | 0.197 | 0.155 | 0.135 | 0.113 | 0.089 | 0.069 |
| DietNeRF [18] | 11.85 | 20.63 | 23.83 | 0.633 | 0.778 | 0.823 | 0.314 | 0.201 | 0.173 | 0.243 | 0.101 | 0.068 |
| RegNerf [45] | 18.89 | 22.2 | 24.93 | 0.745 | *0.841* | *0.884* | 0.19 | 0.117 | 0.089 | 0.112 | 0.071 | 0.047 |
| FreeNerf [84] | <u>19.92</u> | <u>23.25</u> | <u>25.38</u> | <u>0.787</u> | <u>0.844</u> | <u>0.888</u> | *0.182* | 0.137 | 0.096 | <u>0.098</u> | 0.068 | 0.046 |
| MixNerf [65] | 18.95 | 22.3 | 25.03 | 0.744 | 0.835 | 0.879 | 0.203 | *0.102* | *0.065* | 0.113 | *0.066* | *0.042* |
| FlipNerf [64] | *19.55* | *22.45* | 25.12 | 0.767 | 0.839 | 0.882 | <u>0.18</u> | <u>0.098</u> | <u>0.062</u> | *0.101* | *0.064* | *0.041* |
| DiffusioNerf [79] | 16.2 | 20.34 | *25.18* | 0.698 | 0.818 | 0.883 | 0.207 | 0.139 | 0.095 | 0.146 | 0.081 | 0.047 |
| Ours (SparseCraft) | **20.55** | **23.72** | **26.03** | **0.832** | **0.888** | **0.917** | **0.116** | **0.074** | **0.058** | **0.084** | **0.052** | **0.037** |

**Table 2:** Quantitative comparison on DTU. We present the PSNR, SSIM, VGG LPIPS and Average scores of foreground objects. Best scores are in **bold**, second best are <u>underlined</u> and third best are in *italic*.

reconstruction are not the same). Each scene is evaluated on two sets of 3 different views. We use Chamfer distance as metric and report the average of the two sets for each scene, similarly to [40, 58]. We use the same evaluation script as this benchmark, *i.e.* cleaning the generated meshes with masks of training views and sampling points from the generated meshes. We further test our method on few challenging scenes from BlendedMVS [85]. We also evaluate our method on large-scale scenes from Tanks & Temples dataset [27] using only 24 views from the total of more than 150.

For DTU, we report the evaluation as in [37, 58]. We compare mainly against the previously introduced conditional models [4, 88], generalizable reconstruction methods [37, 40, 58] as well as per-scene optimization based neural surface reconstruction methods [46, 72, 86, 87] and the fine-tuned SparseNeuS [40] (denoted SparseNeuS ft). We note that the generalizable methods VolRecon [58] and ReTR [37] do not allow per-scene fine-tuning. We also compare against COLMAP [62].

As there is no standard benchmark for BlendedMVS [85], we use it for qualitative evaluation. We compare against MonoSDF [89], that uses monocular depth and normal priors, COLMAP [62], NeuralAngelo [35] the state-of-the-art hash-based reconstruction method in the dense setting, and S-VolSDF [77], a method that improves the performance of deep MVS through the volumetric rendering of VolSDF [86].

We also compare our method qualitatively on the large-scale dataset of Tanks & Temples [27] against data prior based and test-time optimization method S-VolSDF [77], and NeuralAngelo [35].

## 5.2   Surface Reconstruction

As shown in Table 1, our method SparseCraft outperforms the SOTA on average and on most scenes, even against the generalizable models that were pretrained on other scenes of the same datasets. In particular, we show substantial improvement

for challenging scenes with shiny objects such as scans 110 and 37 as shown in
the qualitative comparison. This showcases also that our method improves largely
on the leveraged result from MVS (COLMAP [62]), as the latter is known for
struggling with shiny/reflective surfaces.

As for qualitative visualizations on both DTU [20] and BlendedMVS [85]
presented in Figures 2 and 3, our method generates overall more detailed and
complete surfaces compared to previous methods. For instance, for scan 118
from DTU [20], NeuS [72] generates an inaccurate surface with many wholes.
The fine-tuned SparseNeuS [40] generates a relatively complete surface, but
overly smooth and lacking important details. VolRecon [58] displays more details
compared to SparseNeuS-ft, but the surface normals appear to be noisy and
inaccurate. Our method can achieve such performance while being faster than
the per-scene optimization methods (Table 6). We note that obtaining the MVS
point cloud (COLMAP) takes only 41 seconds in 3 views, 180 seconds for 9 views
in DTU. Our surfaces show more fidelity and a better trade-off between details
and smoothness. On the challenging large scale dataset T&T (Figure 4), we found
that the generalizable VolRecon [58] fails to generate reasonable outputs. Notice
that our reconstructions display more fidelity and details and fewer failures in
this large scale setting, even-though only a limited number of views is used.

### 5.3   Novel View Synthesis

As shown in table 2, for all input setting, we outperform the current SOTA by
a large margin in all metrics, especially in the most extreme case of 3 input
views. In fact, our method shows superior results on VGG LPIPS which is
reflected in the qualitative comparison in the 3 and 6 input-view settings 5 and 6
respectively, where our renderings appear to be more photo-realistic compared
to RegNeRF [45] and FreeNeRF [84]. For example, the red ball and shiny golden
rabbit show how our method handles well challenging light reflections, and the
example of colorful fruits shows how our method can handle high dynamic range,
all thanks to the various considerations in the design of our method, as well as
the proposed regularization for the sparse setting.

## 6   Ablation

For all reconstruction experiments henceforth, we show average performance over
one of the two sets of 3 views of all 15 DTU scenes. For novel view synthesis
experiments, we report average metrics on all 15 scenes of DTU. We note
that besides the ablations presented below, additional ablations are provided
in the supplementary, including the influence of MVS point cloud's density, the
performance when varying the number of input views, several studies on design
choices related to the sampling of Taylor query points and the effect of progressive
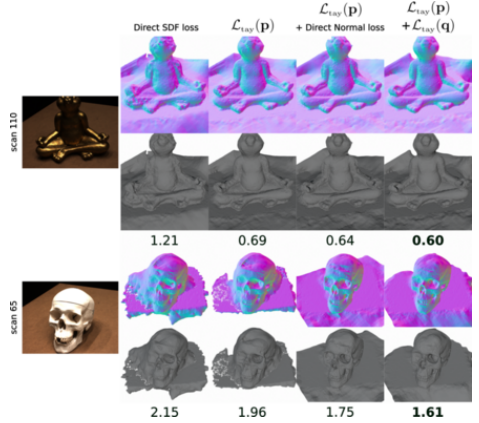hash encoding scheduling.

**Taylor losses _vs_. Direct MVS supervision** We compare our Taylor expansion
based geometric losses to their baselines. The input Taylor loss $\mathcal{L}_{\text{tay}}(\mathbf{p})$ baseline

| Applied losses | Chamfer Distance ↓ |
|---|---|
| SDF loss | 1.34 |
| $\mathcal{L}_{\text{tay}}(\mathbf{p})$ | 1.17 |
| $\mathcal{L}_{\text{tay}}(\mathbf{p})$ + Normal loss | 1.10 |
| $\mathcal{L}_{\text{tay}}(\mathbf{p})$ + $\mathcal{L}_{\text{tay}}(\mathbf{q})$ | **1.08** |

**Table 3:** Numerical Ablation of our Taylor based geometric regularization losses.

| All losses | × | | × | |
|---|---|---|---|---|
| Prog. Enc. | | | × | × |
| Chamfer ↓ | 4.11 | 1.16 | 2.56 | **1.01** |
| PSNR ↑ | | 15.65 | 16.14 | 18.06 | **20.55** |

**Table 4:** Numerical ablation of progressive hash encoding, and the MVS based regularization losses.
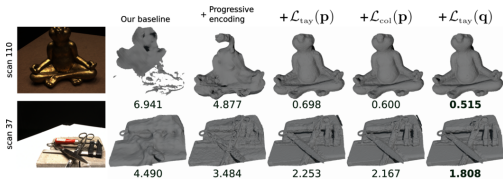


**Fig. 7:** Ablation of our Taylor based geometric regularization losses. We report the **Chamfer score** of reconstructions.

is direct SDF zero supervision. The query Taylor loss $\mathcal{L}_{\text{tay}}(\mathbf{q})$ baseline is a direct SDF gradient supervision with the normal $\mathbf{n}_{\text{MVS}}(\mathbf{p})$. As can be seen in Figure 7 and Table 3, our proposed losses outperform the other combinations and displays improved details and less noise in the reconstructions. This can also be witnessed by the chamfer scores reported in Figure 7. While direct MVS supervision, especially the normal loss, improves our baseline, incorporating this supervision through our Taylor losses is more beneficial. We also find that the proposed Input Taylor loss is largely superior to applying only the SDF zero supervision, which validate our hypothesis about the benefit of enforcing linearity close to surface points.

**Ablation of progressive hash encoding** Table 4 shows improvement brought by the progressive hash encoding, and the MVS based regularization losses. We find that while our regularization can improve surface quality, using progressive encoding act as a regularization and helps to avoid artifacts in the reconstruction, so that the geometry model does not prematurely overfit to fine details. In addition, as our proposed losses are geometric in nature, they sacrifice rendering quality at the expense of good reconstruction. Combining them with the progressive encoding leads to superior rendering quality than using only the progressive encoding.
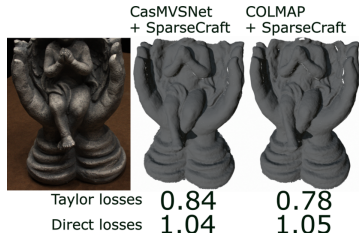
**Ablation of regularization losses** Table 5 and Figure 5 illustrate the contribution of each of our regularization losses to our final performance for both reconstruction and rendering quality. Our baseline model in this case is our method without progressive encoding and MVS regularization. We find that both Taylor-based losses are crucial for learning good surfaces. Further, regularizing the diffuse component of the color network, with MVS color labels alleviates the issue of bias found in the rendering process of NeuS as studied in [10], and hence improves the performance as well while enhancing rendering results even further.

**Fig. 8:** Ablation of our method's components. We report the **Chamfer score** of reconstructions.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $\mathcal{L}_{\mathrm{tay}}(\mathbf{p})$ | × | × | × | × | | | |
| $\mathcal{L}_{\mathrm{tay}}(\mathbf{q})$ | × | × | | | × | | |
| $\mathcal{L}_{\mathrm{col}}(\mathbf{p})$ | × | | × | | × | | |
| Chamfer ↓ | **1.01** | 1.08 | 1.09 | 1.17 | 1.14 | 2.41 | 4.11 |
| PSNR ↑ | **20.55** | 19.65 | 19.44 | 19.27 | 20.12 | 18.05 | 15.65 |
| LPIPS ↓ | **0.116** | 0.132 | 0.146 | 0.152 | 0.127 | 0.192 | 0.269 |

**Table 5:** Numerical ablation of our MVS based regularization losses.



**Fig. 9:** Running our method with learnable MVS. We report the **Chamfer score** of reconstruction.

| Method | Training Time ↓ |
|---|---|
| Ours (SparseCraft) | 9 minutes |
| Ours (SparseCraft) w/o reg. | 7 minutes |
| SparseNeuS ft [40] | 20 minutes |
| NeuralAngelo [35] | 15 minutes |
| S-VolSDF [77] | 18 minutes |
| MonoSDF [89] | 1.5 hours+ |

**Table 6:** Training time on an NVIDIA RTX A6000 of per-scene optimization methods for surface reconstruction from 3 views.

**Using learnable MVS** Figure 9 shows the compatibility of our method with other Point Clouds sources, in this case CasMVSNet [14]. Notice that our novel Taylor losses improve over standard direct losses both when using COLMAP and CasMVSNet.

## 7    Limitations

Since our method uses MVS cues, it suffers from the same limitations of the used MVS method (COLMAP in our case); Thus, it requires enough overlap between input images, and it may not be suitable for reconstruction of highly non-Lambertian surfaces, for which COLMAP is known to fail. In addition, The point cloud obtained from the MVS method have to be dense enough for more accurate normals estimation used by our method. We showed in our experiments how these limitations could be alleviated to some extent by using more advanced MVS techniques such as learnable MVS.

## 8    Conclusion

We presented a new method called SparseCraft for time efficient learning of SDF and radiance fields from sparse imagery. We bridged photogrammetry and deep learning based INRs through novel regularization losses to obtain the SOTA in novel view synthesis and reconstruction simultaneously. Through input data requirement reduction, we hope this work will contribute towards more accessible 3D capture.

# References

1. Aliev, K.A., Sevastopolsky, A., Kolos, M., Ulyanov, D., Lempitsky, V.: Neural point-based graphics. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16. pp. 696–712. Springer (2020)
2. Boulch, A., Marlet, R.: Poco: Point convolution for surface reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6302–6314 (2022)
3. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., De Mello, S., Gallo, O., Guibas, L.J., Tremblay, J., Khamis, S., et al.: Efficient geometry-aware 3d generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16123–16133 (2022)
4. Chen, A., Xu, Z., Zhao, F., Zhang, X., Xiang, F., Yu, J., Su, H.: Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14124–14133 (2021)
5. Chen, C., Han, Z., Liu, Y.S.: Unsupervised inference of signed distance functions from single sparse point clouds without learning priors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
6. Chibane, J., Bansal, A., Lazova, V., Pons-Moll, G.: Stereo radiance fields (srf): Learning view synthesis for sparse views of novel scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7911–7920 (2021)
7. Deng, K., Liu, A., Zhu, J.Y., Ramanan, D.: Depth-supervised nerf: Fewer views and faster training for free. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12882–12891 (2022)
8. Du, Y., Zhang, Y., Yu, H.X., Tenenbaum, J.B., Wu, J.: Neural radiance flow for 4d view synthesis and video processing. In: 2021 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 14304–14314. IEEE Computer Society (2021)
9. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: CVPR (2017)
10. Fu, Q., Xu, Q., Ong, Y.S., Tao, W.: Geo-Neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. CoRR **2205**, 15848 (2022)
11. Furukawa, Y., Ponce, J.: Accurate, dense, and robust multiview stereopsis. IEEE TPAMI **32**(8), 1362–1376 (2009)
12. Galliani, S., Lasinger, K., Schindler, K.: Massively parallel multiview stereopsis by surface normal diffusion. ICCV pp. 873–881 (2015)
13. Gropp, A., Yariv, L., Haim, N., Atzmon, M., Lipman, Y.: Implicit geometric regularization for learning shapes. arXiv preprint arXiv:2002.10099 (2020)
14. Gu, X., Fan, Z., Zhu, S., Dai, Z., Tan, F., Tan, P.: Cascade cost volume for high-resolution multi-view stereo and stereo matching. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2495–2504 (2020)
15. Guo, Y.C.: Instant neural surface reconstruction (2022), https://github.com/bennyguo/instant-nsr-pl
16. Jain, A., Mildenhall, B., Barron, J.T., Abbeel, P., Poole, B.: Zero-shot text-guided object generation with dream fields (2022)
17. Jain, A., Mildenhall, B., Barron, J.T., Abbeel, P., Poole, B.: Zero-shot text-guided object generation with dream fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 867–876 (2022)

18. Jain, A., Tancik, M., Abbeel, P.: Putting nerf on a diet: Semantically consistent few-shot view synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5885–5894 (2021)
19. Jena, S., Multon, F., Boukhayma, A.: Neural mesh-based graphics. In: European Conference on Computer Vision. pp. 739–757. Springer (2022)
20. Jensen, R., Dahl, A., Vogiatzis, G., Tola, E., Aanæs, H.: Large scale multi-view stereopsis evaluation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 406–413 (2014)
21. Johari, M.M., Lepoittevin, Y., Fleuret, F.: Geonerf: Generalizing nerf with geometry priors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18365–18375 (2022)
22. Kajiya, J.T., Von Herzen, B.P.: Ray tracing volume densities. ACM SIGGRAPH computer graphics **18**(3), 165–174 (1984)
23. Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: CVPR (2018)
24. Kazhdan, M., Hoppe, H.: Screened poisson surface reconstruction. ACM Transactions on Graphics (ToG) **32**(3), 1–13 (2013)
25. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics **42**(4), 1–14 (2023)
26. Kim, M., Seo, S., Han, B.: Infonerf: Ray entropy minimization for few-shot neural volume rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 12912–12921 (June 2022)
27. Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and temples: Benchmarking large-scale scene reconstruction. ACM Transactions on Graphics (ToG) **36**(4), 1–13 (2017)
28. Kostrikov, I., Horbert, E., Leibe, B.: Probabilistic labeling cost for high-accuracy multi-view reconstruction. In: Proceedings of the ieee conference on computer vision and pattern recognition. pp. 1534–1541 (2014)
29. Kutulakos, K.N., Seitz, S.M.: A theory of shape by space carving. International journal of computer vision **38**, 199–218 (2000)
30. Lhuillier, M., Quan, L.: A quasi-dense approach to surface reconstruction from uncalibrated images. IEEE transactions on pattern analysis and machine intelligence **27**(3), 418–433 (2005)
31. Li, Q., Multon, F., Boukhayma, A.: Learning generalizable light field networks from few images. In: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1–5. IEEE (2023)
32. Li, Q., Multon, F., Boukhayma, A.: Regularizing neural radiance fields from sparse rgb-d inputs. In: 2023 IEEE International Conference on Image Processing (ICIP). pp. 2320–2324. IEEE (2023)
33. Li, R., Gao, H., Tancik, M., Kanazawa, A.: Nerfacc: Efficient sampling accelerates nerfs. arXiv preprint arXiv:2305.04966 (2023)
34. Li, T., Slavcheva, M., Zollhoefer, M., Green, S., Lassner, C., Kim, C., Schmidt, T., Lovegrove, S., Goesele, M., Newcombe, R., et al.: Neural 3d video synthesis from multi-view video. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5521–5531 (2022)
35. Li, Z., Müller, T., Evans, A., Taylor, R.H., Unberath, M., Liu, M.Y., Lin, C.H.: Neuralangelo: High-fidelity neural surface reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8456–8465 (2023)

36. Li, Z., Niklaus, S., Snavely, N., Wang, O.: Neural scene flow fields for space-time view synthesis of dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6498–6508 (2021)
37. Liang, Y., He, H., Chen, Y.C.: Retr: Modeling rendering via transformer for generalizable neural surface reconstruction. In: Thirty-seventh Conference on Neural Information Processing Systems (2023)
38. Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf: Bundle-adjusting neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5741–5751 (2021)
39. Liu, Y., Peng, S., Liu, L., Wang, Q., Wang, P., Theobalt, C., Zhou, X., Wang, W.: Neural rays for occlusion-aware image-based rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7824–7833 (2022)
40. Long, X., Lin, C., Wang, P., Komura, T., Wang, W.: Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In: Proceedings of the European conference on computer vision (ECCV). pp. 210–227 (2022)
41. Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duckworth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7210–7219 (2021)
42. Max, N.: Optical models for direct volume rendering. IEEE Transactions on Visualization and Computer Graphics $\mathbf{1}$(2), 99–108 (1995)
43. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Nerf, R.N.: Representing scenes as neural radiance fields for view synthesis. ECCV (2020)
44. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. ACM Transactions on Graphics (ToG) $\mathbf{41}$(4), 1–15 (2022)
45. Niemeyer, M., Barron, J.T., Mildenhall, B., Sajjadi, M.S., Geiger, A., Radwan, N.: Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5480–5490 (2022)
46. Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5589–5599 (2021)
47. Ouasfi, A., Boukhayma, A.: Few'zero level set'-shot learning of shape signed distance functions in feature space. In: ECCV (2022)
48. Ouasfi, A., Boukhayma, A.: Few-shot unsupervised implicit neural shape representation learning with spatial adversaries. In: Forty-first International Conference on Machine Learning (2024), https://openreview.net/forum?id=SLqdDWwibH
49. Ouasfi, A., Boukhayma, A.: Mixing-denoising generalizable occupancy networks. In: 3DV (2024)
50. Ouasfi, A., Boukhayma, A.: Robustifying generalizable implicit shape networks with a tunable non-parametric model. Advances in Neural Information Processing Systems $\mathbf{36}$ (2024)
51. Ouasfi, A., Boukhayma, A.: Unsupervised occupancy learning from sparse point cloud. CVPR (2024)
52. Park, K., Sinha, U., Barron, J.T., Bouaziz, S., Goldman, D.B., Seitz, S.M., Martin-Brualla, R.: Nerfies: Deformable neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5865–5874 (2021)

53. Peng, S., Niemeyer, M., Mescheder, L., Pollefeys, M., Geiger, A.: Convolutional occupancy networks. In: Proceedings of the European conference on computer vision (ECCV). pp. 523–540 (2020)
54. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988 (2022)
55. Pumarola, A., Corona, E., Pons-Moll, G., Moreno-Noguer, F.: D-nerf: Neural radiance fields for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10318–10327 (2021)
56. Qin, C., Martens, J., Gowal, S., Krishnan, D., Dvijotham, K., Fawzi, A., De, S., Stanforth, R., Kohli, P.: Adversarial robustness through local linearization. Advances in Neural Information Processing Systems **32** (2019)
57. Rebain, D., Jiang, W., Yazdani, S., Li, K., Yi, K.M., Tagliasacchi, A.: Derf: Decomposed radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14153–14161 (2021)
58. Ren, Y., Wang, F., Zhang, T., Pollefeys, M., Süsstrunk, S.: Volrecon: Volume rendering of signed ray distance functions for generalizable multi-view reconstruction. arXiv preprint arXiv:2212.08067 (2022)
59. Roessle, B., Barron, J.T., Mildenhall, B., Srinivasan, P.P., Nießner, M.: Dense depth priors for neural radiance fields from sparse input views. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12892–12901 (2022)
60. Sara Fridovich-Keil and Alex Yu, Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR (2022)
61. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4104–4113 (2016)
62. Schönberger, J.L., Zheng, E., Pollefeys, M., Frahm, J.M.: Pixelwise view selection for unstructured multi-view stereo. In: European Conference on Computer Vision (ECCV) (2016)
63. Seitz, S.M., Dyer, C.R.: Photorealistic scene reconstruction by voxel coloring. International journal of computer vision **35**, 151–173 (1999)
64. Seo, S., Chang, Y., Kwak, N.: Flipnerf: Flipped reflection rays for few-shot novel view synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 22883–22893 (October 2023)
65. Seo, S., Han, D., Chang, Y., Kwak, N.: Mixnerf: Modeling a ray with mixture density for novel view synthesis from sparse inputs. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 20659–20668 (June 2023)
66. Sitzmann, V., Zollhöfer, M., Wetzstein, G.: Scene representation networks: Continuous 3d-structure-aware neural scene representations. Advances in Neural Information Processing Systems **32** (2019)
67. Srinivas, S., Matoba, K., Lakkaraju, H., Fleuret, F.: Efficient training of low-curvature neural networks. Advances in Neural Information Processing Systems **35**, 25951–25964 (2022)
68. Sun, C., Sun, M., Chen, H.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: CVPR (2022)
69. Trevithick, A., Yang, B.: Grf: Learning a general radiance field for 3d representation and rendering. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15182–15192 (2021)

70. Wang, J., Wang, P., Long, X., Theobalt, C., Komura, T., Liu, L., Wang, W.: Neuris: Neural reconstruction of indoor scenes using normal priors. In: European Conference on Computer Vision. pp. 139–155. Springer (2022)

71. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single rgb images. In: ECCV (2018)

72. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689 (2021)

73. Wang, Q., Wang, Z., Genova, K., Srinivasan, P.P., Zhou, H., Barron, J.T., Martin-Brualla, R., Snavely, N., Funkhouser, T.: Ibrnet: Learning multi-view image-based rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4690–4699 (2021)

74. Wang, Y., Han, Q., Habermann, M., Daniilidis, K., Theobalt, C., Liu, L.: Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 3295–3306 (October 2023)

75. Wang, Y., Skorokhodov, I., Wonka, P.: Hf-neus: Improved surface reconstruction using high-frequency details. Advances in Neural Information Processing Systems **35**, 1966–1978 (2022)

76. Wei, Y., Liu, S., Rao, Y., Zhao, W., Lu, J., Zhou, J.: Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5610–5619 (2021)

77. Wu, H., Graikos, A., Samaras, D.: S-volsdf: Sparse multi-view stereo regularization of neural implicit surfaces. arXiv preprint arXiv:2303.17712 (2023)

78. Wu, T., Wang, J., Pan, X., Xu, X., Theobalt, C., Liu, Z., Lin, D.: Voxurf: Voxel-based efficient and accurate neural surface reconstruction (2022). `https://doi.org/10.48550/ARXIV.2208.12697`, `https://arxiv.org/abs/2208.12697`

79. Wynn, J., Turmukhambetov, D.: Diffusionerf: Regularizing neural radiance fields with denoising diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4180–4189 (2023)

80. Xian, W., Huang, J.B., Kopf, J., Kim, C.: Space-time neural irradiance fields for free-viewpoint video. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9421–9431 (2021)

81. Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., Neumann, U.: Point-nerf: Point-based neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5438–5448 (2022)

82. Xu, Q., Tao, W.: Multi-scale geometric consistency guided multi-view stereo. CVPR pp. 5483–5492 (2019)

83. Yang, B., Bao, C., Zeng, J., Bao, H., Zhang, Y., Cui, Z., Zhang, G.: Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In: European Conference on Computer Vision. pp. 597–614. Springer (2022)

84. Yang, J., Pavone, M., Wang, Y.: Freenerf: Improving few-shot neural rendering with free frequency regularization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8254–8263 (2023)

85. Yao, Y., Luo, Z., Li, S., Zhang, J., Ren, Y., Zhou, L., Fang, T., Quan, L.: Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1790–1799 (2020)

86. Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. Advances in Neural Information Processing Systems **34**, 4805–4815 (2021)

87. Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.:
    Multiview neural surface reconstruction by disentangling geometry and appearance.
    Advances in Neural Information Processing Systems **33**, 2492–2502 (2020)
88. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from
    one or few images. In: Proceedings of the IEEE/CVF Conference on Computer
    Vision and Pattern Recognition. pp. 4578–4587 (2021)
89. Yu, Z., Peng, S., Niemeyer, M., Sattler, T., Geiger, A.: Monosdf: Exploring monoc-
    ular geometric cues for neural implicit surface reconstruction. Advances in neural
    information processing systems **35**, 25018–25032 (2022)
90. Zhang, J., Yang, G., Tulsiani, S., Ramanan, D.: Ners: Neural reflectance surfaces
    for sparse-view 3d reconstruction in the wild. Advances in Neural Information
    Processing Systems **34**, 29835–29847 (2021)
91. Zhang, J., Yao, Y., Li, S., Fang, T., McKinnon, D., Tsin, Y., Quan, L.: Critical
    regularizations for neural surface reconstruction in the wild. In: Proceedings of
    the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp.
    6270–6279 (2022)
92. Zhang, Y., Hu, Z., Wu, H., Zhao, M., Li, L., Zou, Z., Fan, C.: Towards unbiased
    volume rendering of neural implicit surfaces with geometry priors. In: Proceedings
    of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp.
    4359–4368 (2023)