

# SparseRecon: Neural Implicit Surface Reconstruction from Sparse Views with Feature and Depth Consistencies

Anonymous ICCV submission

Paper ID 2711

## Abstract

001 Surface reconstruction from sparse views aims to recon-  
 002 struct a 3D shape or scene from few RGB images. How-  
 003 ever, existing generalization-based methods do not gen-  
 004 eralize well on views that were unseen during training,  
 005 while the reconstruction quality of overfitting-based meth-  
 006 ods is still limited by the limited geometry clues. To ad-  
 007 dress this issue, we propose SparseRecon, a novel neural  
 008 implicit reconstruction method for sparse views with vol-  
 009 ume rendering-based feature consistency and uncertainty-  
 010 guided depth constraint. Firstly, we introduce a feature con-  
 011 sistency loss across views to constrain the neural implicit  
 012 field. This design alleviates the ambiguity caused by insuf-  
 013 ficient consistency information of views and ensures com-  
 014 pleteness and smoothness in the reconstruction results. Sec-  
 015 ondly, we employ an uncertainty-guided depth constraint  
 016 to back up the feature consistency loss in areas with oc-  
 017 clusion and insignificant features, which recovers geometry  
 018 details for better reconstruction quality. Experimental re-  
 019 sults demonstrate that our method outperforms the state-of-  
 020 the-art methods, which can produce high-quality geometry  
 021 with sparse-view input, especially in the scenarios on small  
 022 overlapping views.

## 023 1. Introduction

024 As one of the important tasks in computer vision, 3D re-  
 025 construction has attracted lots of research attentions in re-  
 026 cent years. With the advancement of deep learning, 3D re-  
 027 construction using neural implicit representations based on  
 028 point clouds [25, 32, 54, 55] or images [22, 35, 44, 56] be-  
 029 comes a popular research topic. Although existing methods  
 030 [5, 35, 37, 40, 44, 52] that directly use images have made  
 031 great progress in terms of the reconstruction quality and re-  
 032 construction speed, they require a large number of dense  
 033 views as supervision. When the number of available views  
 034 is limited, current reconstruction methods usually struggle  
 035 to reconstruct high-quality surfaces.

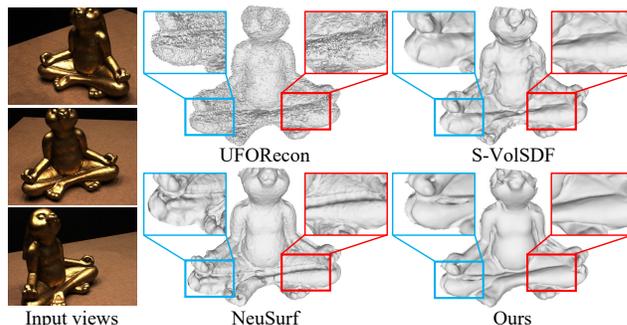


Figure 1. Given only 3 input images with large view angle change, our method can reconstruct a smoother surface compared to the state-of-the-art methods, such as UFORecon [27], S-VolSDF [39] and NeuSurf [13]. The details of each surface are shown in the colored boxes.

Existing methods for sparse view reconstruction can be mainly classified into two categories: generalization-based methods and overfitting-based methods. The generalization-based methods [21, 23, 27, 29, 30] emphasize the generalization of sparse-view reconstruction, but they are mainly effective in scenarios with large view overlaps. In cases with views that were unseen during training, the quality of the reconstructed surface degenerates significantly, as shown in Figure 1. Meanwhile, it takes a long time to pre-train these methods on large-scale data. Instead, overfitting-based methods [13, 14, 39, 45, 46] typically fit the 3D geometry directly from the sparse views by leveraging geometry clues. They show promising capability in reconstructing higher-quality geometric surfaces with small-overlapping views. However, the reconstruction quality of the existing methods is still unsatisfactory.

In this paper, we introduce a multi-view feature consistency loss based on volume rendering and an uncertainty-guided depth constraint to learn neural signed distance functions. This approach allows us to achieve high-quality mesh reconstruction on more challenging sparse views with small overlap.

058 For the *feature consistency loss*, we first employ the pre-  
 059 trained Vis-MVSNet [49] to obtain depth features from the  
 060 input images. Then, within a neural implicit rendering  
 061 framework, the sampled 3D points along the rays emitted  
 062 from the reference image are projected to the source image  
 063 and the reference image. This allows us to acquire source  
 064 features and reference features of each 3D point and mea-  
 065 sure the similarity between these two kinds of features. Fi-  
 066 nally, the feature similarity for each 3D point along the rays  
 067 is accumulated through volume rendering, thus yielding the  
 068 feature similarity associated with the rays. During opti-  
 069 mization, we pursue higher feature similarity along the rays.  
 070 Since the depth information is implicitly encoded with image  
 071 features, feature consistency constraint can significantly  
 072 alleviate the ambiguity issues arising from insufficient con-  
 073 sistency of sparse views and low-texture during reconstruc-  
 074 tion.

075 For the *uncertainty-guided depth prior constraint*, we  
 076 follow MonoSDF [46], utilizing a pre-trained network to  
 077 acquire depth priors for each image, and then use it to con-  
 078 strain the regions with uncertain depth. However, monocular  
 079 depth priors do not have consistent scales to the ground  
 080 truth depth, which are hard to get calibrated to ground truth  
 081 either due to the distortion. To effectively leverage the  
 082 depth priors and provide proper supervision for occluded  
 083 or under-constrained regions, we propose an uncertainty-  
 084 guided depth prior constraint. First, we calibrate the depth  
 085 priors using sparse point clouds obtained from COLMAP  
 086 [31]. Then, during training, we compute the depth confi-  
 087 dence from the rendered depth and impose the depth prior  
 088 constraint only in regions with low confidence. This con-  
 089 straint helps infer more accurate geometry in occluded or  
 090 under-constrained regions, minimizing the negative impact  
 091 of depth prior errors on well-constrained regions.

092 We evaluate our methods on several widely used bench-  
 093 marks and report the state-of-the-art results. In summary,  
 094 our main contributions are as follows.

- 095 • We propose a novel feature consistency loss based on vol-  
 096 ume rendering. It can effectively constrains the neural  
 097 radiance field by leveraging feature consistency among  
 098 multiple views, improving the performance in sparse-  
 099 view reconstruction tasks.
- 100 • By incorporating depth confidence, we utilize the cali-  
 101 brated depth prior more effectively to enhance geometric  
 102 constraints, further improving the reconstruction quality.
- 103 • Extensive experiments on the well-known datasets, such  
 104 as DTU [16] and BlendedMVS[43], demonstrate that our  
 105 method outperforms existing sparse-view reconstruction  
 106 methods and achieve the state-of-the-art results.

## 2. Related Work 107

### 2.1. Neural Implicit Reconstruction 108

Neural implicit reconstruction methods [5, 7, 20, 35–37, 40, 109  
 44, 46], have been rapidly developed based on neural vol- 110  
 ume rendering [26]. These methods introduce the Signed 111  
 Distance Function (SDF) as the implicit representation of 112  
 3D surfaces in volume rendering, achieving multi-view 3D 113  
 reconstruction. While these methods have made significant 114  
 improvements in both reconstruction quality and speed, it 115  
 is important to note that they heavily rely on multiple view- 116  
 points during the optimization. 117

**Generalization-based surface reconstruction with 118  
 sparse views.** In order to directly generalize the reconstruc- 119  
 tion results on sparse views, methods [21, 23, 27, 29, 30, 41] 120  
 adopt the strategy of aggregating features from multiple 121  
 view images to construct a feature volume, which is then 122  
 used to predict the SDF for reconstructing the surface. Vol- 123  
 Recon [30] uses transformers [17] to aggregate multi-view 124  
 features, C2F2NeUS [41] employs cascade architecture to 125  
 construct a volume pyramid, while ReTR [21] and UFORe- 126  
 con [27] aggregates multi-level features. These methods 127  
 require pretraining on large-scale datasets, which typically 128  
 takes several days. Moreover, when there is a significant 129  
 domain gap between the testing and training data, they all 130  
 fail to reconstruct shapes effectively. 131

**Overfitting-based surface reconstruction with sparse 132  
 views.** In contrast, overfitting-based methods directly fit 133  
 the 3D geometry from the sparse images by geometric prior 134  
 constraints. MonoSDF [46] employs depth and normal pri- 135  
 ors to achieve sparse reconstruction with small-overlapping 136  
 views. However, such priors come with errors, and it does 137  
 not fully leverage inter-view consistency, resulting in lower 138  
 reconstruction quality. S-VolSDF [39] employs probability 139  
 volumes obtained from MVS [9] models to guide the ren- 140  
 dering weight estimated by VolSDF [44]. This improves the 141  
 reconstruction results in sparse views with small overlap. 142  
 However, the uncertainties in volumes make negative im- 143  
 pact on the reconstruction surface, leading to surface rough- 144  
 ness or significant defects. More recently, NeuSurf [13] 145  
 leverages sparse point clouds and employs CAP-UDF [54] 146  
 to construct an implicit geometric prior to improve the re- 147  
 construction quality of sparse views. However, when the 148  
 sparse point cloud fails to cover the majority of positions on 149  
 the object surface, effective implicit geometric prior infor- 150  
 mation cannot be obtained, which does not improve the re- 151  
 construction quality. In contrast, our method employs more 152  
 robust feature priors, calculates feature consistency based 153  
 on volume rendering, and simultaneously utilizes depth pri- 154  
 ors to optimize the occluded regions, ultimately resulting in 155  
 high-quality geometric surface. 156

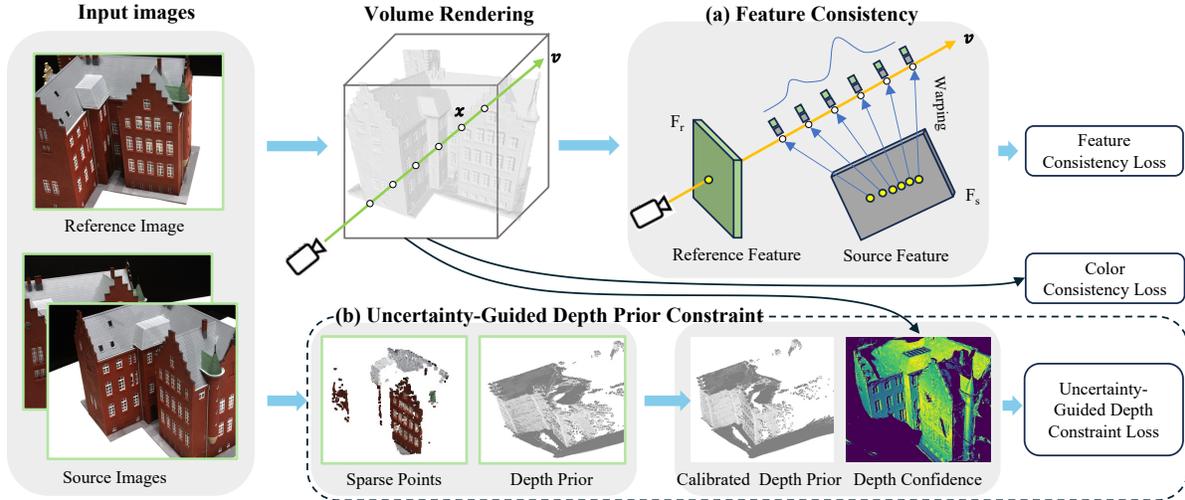


Figure 2. SparseRecon consists of two main parts. (a) Volume rendering-based feature consistency constraint. We extract features from the reference image and source images. For a ray emitted from the reference image, we project each sampled point on the ray onto the source images to obtain the corresponding features. Then, the volume rendering-based feature consistency loss is calculated using the corresponding features on the reference image. (b) Uncertainty-guided depth prior constraint. We use another pre-trained network to obtain the depth prior of the reference image and calibrate it with the sparse point cloud obtained by COLMAP. Then, we calculate the confidence of the rendered depth, so that the calibrated depth prior only constrains areas with low confidence.

157 **2.2. Gaussian Splatting.**

158 Gaussian Splatting [18] has achieved unprecedented opti-  
 159 mization speed and rendering quality in the task of novel  
 160 view synthesis. However, since the Gaussians are unorga-  
 161 nized, the discrete and unstructured points make it difficult  
 162 to extract 3D surfaces through post-processing. To address  
 163 this issue, some methods introduce regularization terms  
 164 [10], convert 3D Gaussians to 2D surfels [4, 12], acquire  
 165 opacity fields through rays [47], improve the depth render-  
 166 ing algorithm [1] of 3DGS, or jointly optimize 3DGS with  
 167 neural radiance fields [2, 24, 53]. However, these meth-  
 168 ods are only applicable to dense views. Recently, FatesGS  
 169 [14] achieves fast sparse-view reconstruction by leveraging  
 170 depth priors and on-surface feature consistency constraints.  
 171 However, due to the poor convergence of the on-surface fea-  
 172 ture consistency constraint and the inaccuracy of the depth  
 173 priors, the reconstruction results still exhibit roughness or  
 174 noticeable defects.

175 **2.3. Sparse View Synthesis.**

176 In addition, the novel view synthesis from sparse views  
 177 is another category of work closely related to sparse view  
 178 reconstruction. Depending on the technical framework,  
 179 these works can be categorized into NeRF-based methods  
 180 [6, 15, 28, 33, 34, 42, 48] and Gaussian Splatting-based  
 181 methods [3, 11, 19, 51, 57]. This line of research also em-  
 182 ploys a limited number of views as input. However, they  
 183 solely focus on the rendering quality of novel views rather  
 184 than surface reconstruction, which are not designed specifi-

cally for the accurate geometric surface reconstruction. Due  
 to the discernible bias (i.e. inherent geometric errors) [35]  
 caused by the conventional volume rendering method or in-  
 consistencies in depth that appear in Gaussian rendering,  
 current sparse view synthesis methods still fail to correctly  
 reconstruct high-fidelity geometric surfaces.

**3. Method**

The overview of our method is depicted in Figure 2. We in-  
 troduce a novel feature consistency loss and an uncertainty-  
 guided depth constraint based on the NeuS [35] framework.  
 In this section, we first explain how to compute feature con-  
 sistency for sampled points along rays. Then we explain  
 how to enhance geometric constraints using depth priors  
 and depth uncertainty. Thirdly, we introduce the color con-  
 sistency loss. Finally, we present the overall loss function  
 for optimization.

**3.1. Volume Rendering-based Feature Consistency**

First, we use a pre-trained MVS network [49] to extract the  
 features from both the reference image and the source im-  
 age. Given a ray emitted from the reference image, let  $p_r(0)$   
 denote the point where a ray intersects the reference image.  
 And for each point  $x_i$  along the ray, we denote its projection  
 on the source image as  $p_s(i)$ . Then, we bilinearly interpo-  
 late  $F_r(0)$  and  $F_s(i)$  at points  $p_r(0)$  and  $p_s(i)$  on im-  
 age features, respectively. Formally, we define the feature con-

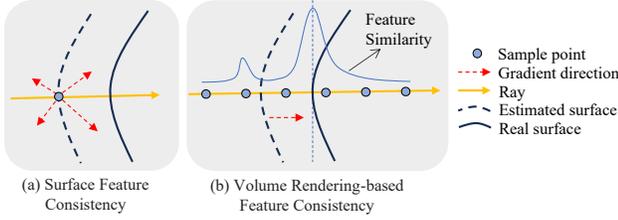


Figure 3. The illustration of (a) on-surface feature consistency and (b) feature consistency with volume rendering.

210 consistency loss function as follows,

$$L_{feat} = M^{occ} \left( 1 - \frac{1}{N} \sum_{i=1}^N w_i f_{cos}(F_r(p_r(0)), F_s(p_s(i))) \right), \quad (1)$$

211 where  $f_{cos}$  is the cosine similarity, and  $w_i$  corresponds to  
 212 the weight for each point along the ray.  $p_s(i) = K(Rx_i + t)$   
 213 is the projection of  $x_i$  in source view, and  $[K; R; t]$  is  
 214 the camera parameters of source view.  $M^{occ}$  is the occlusion  
 215 mask.  
 216

217 Although MVSDf [50], NeuSurf [13] and FatesGS [14]  
 218 also employ feature consistency constraints, they just lever-  
 219 age the intersection point between a camera ray and the ob-  
 220 ject’s surface. Then, this intersection point gets projected  
 221 onto adjacent views to obtain the corresponding image fea-  
 222 tures for the purpose of comparing features at this point  
 223 across multiple views. In sparse view scenarios, the esti-  
 224 mated positions of surface points can easily deviate signifi-  
 225 cantly, making the on-surface feature consistency loss not  
 226 converge. NeuSurf [13] and FatesGS [14] utilize sparse  
 227 point clouds generated by COLMAP [31] as priors, en-  
 228 abling it to obtain partially accurate positions of surface  
 229 points, thereby allowing the on-surface feature consistency  
 230 loss to be more effectively leveraged. However, in regions  
 231 of lacking surface points, the on-surface feature consistency  
 232 loss cannot ensure the attainment of high-quality geometric  
 233 surfaces.

234 Figure 3 illustrates the difference between on-surface  
 235 feature consistency and volume rendering-based feature  
 236 consistency. Due to the uncertainty of gradient direction,  
 237 the constraint solely relying on surface point features is  
 238 challenging to be optimized. In contrast, our method does  
 239 not require the prior estimation of surface points, it calcu-  
 240 lates feature consistency on all sampling points along the  
 241 ray, and provides more reasonable and comprehensive su-  
 242 pervision to the implicit field, thereby addressing the con-  
 243 vergence issue that may arise in sparse reconstruction for  
 244 MVSDf [50] and NeuSurf [13].

### 245 3.2. Uncertainty-Guided Depth Constraint

246 Although multi-view features offer more robust constraint  
 247 than image colors, they are ineffective for occluded regions.

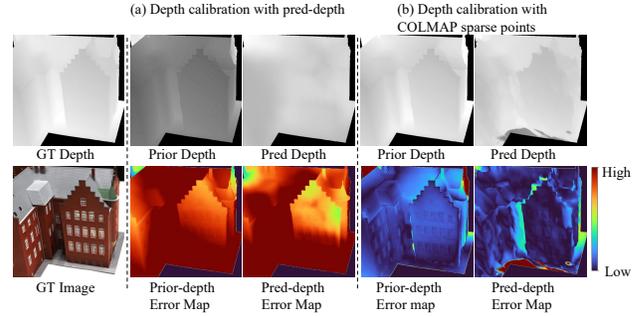


Figure 4. The illustration of predicted depth produced by differ-  
 ent depth prior utilization methods, along with the corresponding  
 error maps. (a) Calibrate the depth prior using the predicted depth  
 during training. (b) Calibrate the depth prior using the COLMAP  
 sparse point cloud.

248 Due to the limited number of views, some regions may only  
 249 be visible from a single viewpoint. To enhance geometric  
 250 constraints, we employ depth priors to supervise the radiance  
 251 field. However, monocular depth priors are not perfect  
 252 and accurate. Although MonoSDF [46] has already takes  
 253 the inaccuracy of depth priors into account, i.e., it aligns  
 254 depth priors using rendered depth during training. How-  
 255 ever, the rendered depth during training is inaccurate, result-  
 256 ing in significant errors in the calibrated depth priors. This  
 257 ultimately leads to the accumulation of errors during train-  
 258 ing, which results in inaccurate reconstructions. Figure 4  
 259 (a) shows the calibrated depth prior and rendered depth ob-  
 260 tained by MonoSDF [46], as well as their error maps com-  
 261 pared to the ground truth depth. It can be seen that both the  
 262 calibrated depth prior and the rendered depth are with large  
 263 errors. Therefore, MonoSDF [46] uses a weight annealing  
 264 strategy to anneal the weight of depth loss to 0 during the  
 265 first 200 training epochs.

266 Another trivial approach is to calibrate the depth priors  
 267 using the sparse point cloud obtained from COLMAP [31].  
 268 Since the sparse points are generally located on the geo-  
 269 metric surface of the object, their depth is relatively accu-  
 270 rate. Therefore, calibrating the depth priors using the sparse  
 271 point cloud can lead to more accurate depth priors. Figure  
 272 4 (b) shows the depth priors calibrated with the sparse point  
 273 cloud, and the depth rendered with the depth priors as a con-  
 274 straint, as well as their error maps compared to the ground  
 275 truth depth. It indicates that the depth priors calibrated to  
 276 the point cloud from the COLMAP [31] are more accurate.  
 277 Therefore, we can use them as an constraint leads to more  
 278 precise rendered depth.

279 However, due to the distortions in monocular depth pri-  
 280 ors, it is impossible to perfectly align them with the ground  
 281 truth depth. Even after calibration, the depth priors still ex-  
 282 hibit noticeable errors when compared to the ground truth  
 283 depth. In sparse view scenarios, occlusions and insufficient

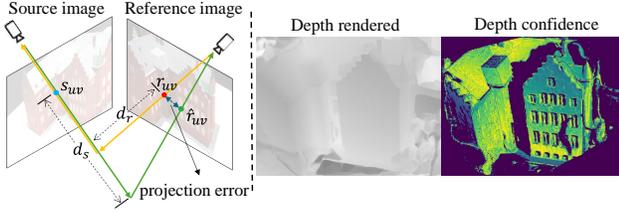


Figure 5. Left: the method of obtaining the confidence of rendered depth. Right: the rendered depth and the depth confidence.

constraints are more common, leading to significant discrepancies between the geometry of occluded regions and the real surface. Therefore, to achieve more accurate geometry in these under-constrained regions while avoiding the negative impact of depth prior errors on well-constrained regions, we propose an uncertainty-guided depth prior constraint method to more effectively utilize the depth priors. Specifically, we apply depth prior constraints in regions with depth uncertainty, while refraining from using them in regions with high depth confidence.

To obtain the confidence of the rendered depth, we employ a method to evaluate the multi-view depth projection consistency. As shown in Figure 5, for a specific pixel  $r_{uv}$  in the reference image with depth  $d_r$ , it can be mapped to a neighboring image through the homography matrix  $H_{rs}$ , leading to a pixel  $s_{uv}$ ,

$$s_{uv} = H_{rs}r_{uv}, \quad (2)$$

$$H_{rs} = M_s M_r^{-1}, \quad (3)$$

where  $M_r$  and  $M_s$  are the projection matrices corresponding to the reference and source views, respectively. Similarly, we can map the pixel  $s_{uv}$  in the source view to the reference view using the projection matrix  $H_{sr}$  and its corresponding depth  $d_s$ , resulting in  $\hat{r}_{uv}$ . The forward and backward projection distance error reflect the accuracy of depth predictions, so we take it as the depth confidence, which is defined as

$$C_d = \begin{cases} \frac{1}{e^{\|r_{uv} - \hat{r}_{uv}\|}}, & \text{if } \|r_{uv} - \hat{r}_{uv}\| \leq 1 \\ 0, & \text{if } \|r_{uv} - \hat{r}_{uv}\| > 1 \end{cases} \quad (4)$$

The right side of Figure 5 shows the rendered depth and the corresponding depth confidence.

Correspondingly, the depth uncertainty is defined as  $U_d = 1 - C_d$ . Meanwhile, we can set a threshold  $\tau$  for depth confidence  $C_d$  to obtain the occlusion mask  $M^{occ} = \{C_d > \tau\}$ .

For depth calibration, we leverage COLMAP [31] to obtain a sparse point cloud  $\{X : x_1, x_2 \dots x_i \in R\}$  and visibility flags indicating which keypoints are visible from view  $I$ . Given the camera parameters  $P$  of view  $I$ , we estimate

the depth  $\bar{D}_i$  of keypoints by computing the distance from the visible keypoints  $x_i$  to the camera center  $o$ . Then, we calibrate the monocular depth prior  $\hat{D}$  with  $\bar{D}_i$ , it can be defined as  $\bar{D} \approx a\hat{D} + b$ , where  $a$  is the scale factor and  $b$  is the shift factor, obtained through the least squares method. Formally, the depth constraint loss is defined as,

$$L_{depth} = \sum_{r \in R} U_d \left\| (a\hat{D} + b) - D_{pred} \right\|^2. \quad (5)$$

### 3.3. Color Consistency Constraint

Although feature consistency constraint can ensure that the reconstruction does not suffer from severe artifacts, it does not provide sufficient supervision to reconstruct fine geometric details. Conversely, in cases with rich textures, image color constraint can refine the geometric details. Therefore, following the NeuralWarp [5], pixel warping loss and patch warping loss are used in our method as multi-view color consistency loss functions,

$$L_{color} = \sum_{r \in R} M^{occ} d_{pixel}(C(r), C_s(r)) + \sum_{r \in R} M^{occ} d_{patch}(P(r), P_s(r)), \quad (6)$$

where  $C(r)$  and  $C_s(r)$  are the ground truth color of the pixel from which the ray emits and the rendered color, respectively,  $P(r)$  and  $P_s(r)$  are the ground truth color of the patch corresponding to the ray and the rendered patch color, respectively.  $d_{pixel}$  is the loss metric for pixel color, where we use  $L1$  loss as  $d_{pixel}$ .  $d_{patch}$  is the loss metric for patch color, where we use the Structural Similarity Measure (SSIM [38]) as  $d_{patch}$ .

### 3.4. Training Loss

In addition to the above-mentioned three loss functions, we also use the Eikonal loss [8] used in NeuS [35]. We define the overall loss function as follows:

$$L = L_{feat} + \alpha L_{depth} + L_{color} + \beta L_{eik}, \quad (7)$$

$L_{eik}$  is the Eikonal loss [8], used to regularize the SDF values of sampled points, defined as

$$L_{eik} = \frac{1}{mn} \sum_{i,k} (\|\nabla f(x_{i,k})\|_2 - 1)^2. \quad (8)$$

## 4. Experiments

### 4.1. Dataset

We evaluate our method on DTU [16] and BlendedMVS [43] dataset. For the DTU [16] dataset, to avoid using the scenes that have already been used as training data on the pretrained Vis-MVSNet [49] model, we select the same 11

Methods	21	24	34	37	38	40	82	106	110	114	118	Mean CD↓
VolSDF [44]	5.47	4.38	3.15	7.38	1.88	6.70	5.19	4.67	2.79	1.32	1.83	4.07
NeuS [35]	5.63	3.58	6.00	4.60	2.57	4.53	1.91	4.18	5.46	1.19	4.16	3.98
NeuralWarp [5]	2.53	1.88	<u>0.74</u>	<u>1.80</u>	<b>0.84</b>	11.50	2.64	2.10	4.37	1.19	2.63	2.93
MonoSDF [46]	4.14	5.92	1.39	4.55	2.19	2.14	2.36	5.62	4.58	1.63	3.02	3.41
Vis-MVSNet [49]	3.39	4.44	0.85	3.36	1.69	3.35	3.35	2.34	2.16	0.74	1.83	2.50
MVSDF [50]	4.31	4.71	1.65	6.37	1.77	4.47	3.61	1.87	1.67	1.25	1.69	3.03
2DGS [12]	4.47	3.54	3.48	4.13	4.25	3.61	4.83	2.40	2.97	1.35	2.17	3.38
PGSR [1]	5.58	4.01	3.15	5.19	4.55	3.65	5.57	2.35	1.91	0.57	1.55	3.46
SparseNeuS <sub>ft</sub> [23]	3.48	4.37	2.92	4.76	2.79	3.73	2.80	1.86	3.10	1.15	2.29	3.02
VolRecon [30]	2.72	3.07	1.82	4.32	2.14	3.04	3.00	2.56	2.81	1.49	3.22	2.75
GenS <sub>ft</sub> [29]	5.86	7.67	3.62	8.57	5.37	5.41	5.48	6.04	5.29	4.69	4.35	5.67
ReTR [21]	2.67	3.37	1.62	3.68	1.87	3.40	3.67	2.84	2.85	1.56	2.35	2.72
UFORecon [27]	<b>1.84</b>	1.52	0.79	2.58	1.00	<u>1.82</u>	<u>1.72</u>	1.20	0.93	0.66	1.26	<u>1.39</u>
S-VolSDF [39]	2.45	3.08	1.33	3.09	1.22	3.21	1.91	1.51	1.23	0.74	1.2	1.91
SparseCraft [45]	2.88	2.42	0.92	2.97	1.58	2.78	2.51	1.10	5.24	0.65	0.88	2.16
NeuSurf [13]	7.60	1.43	2.93	3.18	1.53	2.86	1.86	<u>1.09</u>	1.41	<b>0.37</b>	<b>0.62</b>	2.26
FatesGS [14]	3.98	<u>1.32</u>	2.53	2.85	3.36	2.71	3.76	1.49	<u>0.85</u>	0.47	1.06	2.22
Ours	<u>2.14</u>	<b>1.26</b>	<b>0.72</b>	<b>1.46</b>	<u>0.86</u>	<b>1.39</b>	<b>1.37</b>	<b>0.94</b>	<b>0.77</b>	<u>0.44</u>	<u>0.83</u>	<b>1.11</b>

Table 1. Quantitative results of Chamfer Distance (CD↓) on DTU dataset with 3 *small-overlapping* images. The methods are divided into three categories, from top to bottom: (1) dense-view reconstruction methods related to ours, (2) generalization-based sparse-view reconstruction methods, and (3) overfitting-based sparse-view reconstruction methods. the best results are in *bold*, the second best are *underlined*.

361 scenes as in S-VolSDF [39]. The image resolution is set  
 362 to 1600×1200. Similar to the S-VolSDF [39] and NeuSurf  
 363 [13] methods, we select the views 22, 25, and 28 for the  
 364 more challenging reconstruction of small overlaps.

365 For the BlendedMVS [43] dataset, we follow the S-  
 366 VolSDF [39] to use the same 9 challenging scenes, with 3  
 367 small-overlapping views for each scene. The image resolu-  
 368 tion is set to 768×576.

## 369 4.2. Implementation Details

370 We use the same network architecture and initialization  
 371 strategy as NeuS [35] and incorporated our volume ren-  
 372 dering feature consistency loss, uncertainty-guided depth  
 373 constraint loss, and color consistency loss. For the weight  
 374 factors in the loss functions Eq. 7, we set the  $\alpha$  for the  
 375 uncertainty-guided depth prior constraint loss  $L_{depth}$  to 0.5  
 376 and the  $\beta$  for the Eikonal loss  $L_{eik}$  to 0.1. Each scene  
 377 is trained 100K iterations on a RTX3090 GPU. The patch  
 378 warping term in the color consistency loss requires the sur-  
 379 face point normals to calculate homographies, but the initial  
 380 normals are too noisy [5], therefore, the patch warping loss  
 381 is applied after 20k training steps. The threshold  $\tau$  of the  
 382 occlusion mask is set to 0.

## 383 4.3. Baseline

384 We compare our approach with three categories of meth-  
 385 ods. *Dense-view methods*: NeuS [35], VolSDF [39], Neu-

ralWarp [5], Vis-MVSNet [49], MVSDf [50], 2DGS [12]  
 and PGSR [1]. *Generalization-based methods*: SparseNeuS  
 [23], VolRecon [30], GenS [29], ReTR [21] and UFOre-  
 con [27]. *Overfitting-based methods*: S-VolSDF [39], Spar-  
 seCraft [45], NeuSurf [13] and FatesGS [14]. The recon-  
 struction results for SparseNeuS [23] and GenS [29] are  
 fine-tuned using 3 views for each scene.

## 4.4. Comparisons

**Reconstruction on DTU.** For a comprehensive compar-  
 ison, we evaluate the baselines and our method on both  
 small-overlapping and large-overlapping views. Following  
 baselines [13, 14, 23], we report the Chamfer Distance (CD)  
 between the reconstruction surfaces and the ground truth  
 point clouds. The CD results with small overlapping views  
 are shown in Table 1. The meshes reconstructed by several  
 methods using 3 views with small overlapping are shown  
 in Fig. 6. For the generalization-based sparse reconstruc-  
 tion methods, we only show the reconstruction results of  
 the latest UFORecon [27], as the reconstruction quality of  
 other methods is lower than that of UFORecon [27]. The  
 experimental results show that our method significantly im-  
 proves the mesh quality with small overlap views, compared  
 to the state-of-the-art sparse-view reconstruction methods.  
 The results of large overlapping views are presented in the  
 supplementary materials.

As shown in Figure 6, when input sparse views with

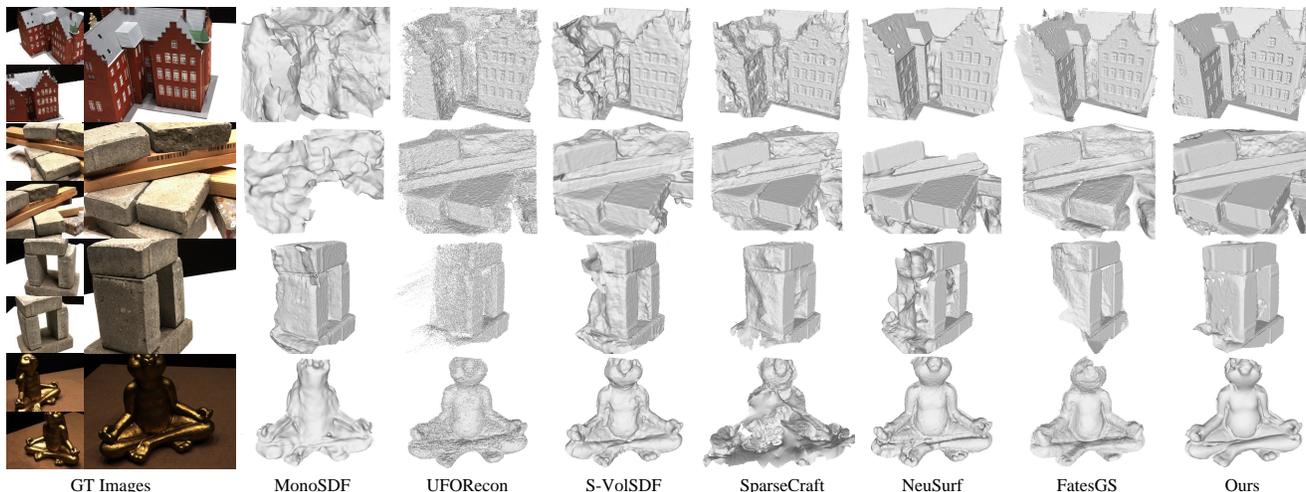


Figure 6. Visual comparison on DTU dataset with 3 *small-overlapping* images.

412 small overlap, both MonoSDF [46] and SparseCraft [45]  
 413 suffer from reconstruction ambiguity and failures, high-  
 414 lighting that relying solely on simplistic geometric prior  
 415 constraints is insufficient to obtain complete and accurate  
 416 meshes. UFORecon [27] shows significant roughness in its  
 417 reconstruction results. S-VolSDF [39], NeuSurf [13] and  
 418 FatesGS [14] exhibit noticeable reconstruction defects. Ex-  
 419 perimental results demonstrate that our method is effective  
 420 in alleviating geometric and appearance ambiguities during  
 421 the optimization process. This significantly enhances the  
 422 quality of mesh reconstruction, especially in scenarios with  
 423 small overlapping views and low texture.

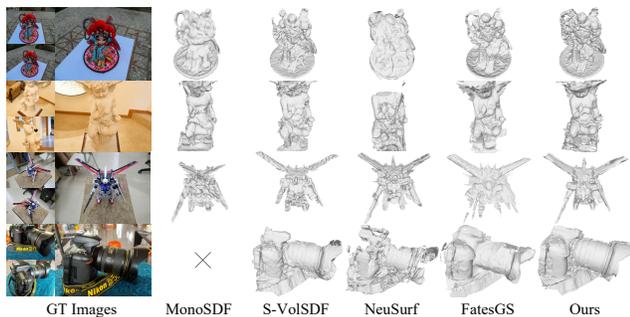


Figure 7. Visual comparison on BlendedMVS dataset. 'x' indicates reconstruction failure.

424 **Reconstruction on BlendedMVS.** Figure 7 presents the  
 425 visual comparison of reconstructed mesh for overfitting-  
 426 based methods. With only 3 small-overlapping views pro-  
 427 vided, all of the generalization-based methods completely  
 428 fail to reconstruct in the sparse setting of BlendedMVS[43]  
 429 dataset, even if SparseNeuS [23] is fine-tuned. Therefore,

430 the reconstruction results of these methods are not included  
 431 in Figure 7. Compared to other methods, our approach can  
 432 generate more complete and detailed meshes. Similarly,  
 433 MonoSDF [46] fails to reconstruct either. The meshes gen-  
 434 erated by S-VolSDF [39], NeuSurf [13] and FatesGS [14]  
 435 exhibit significant defects. Both NeuSurf [13] and FatesGS  
 436 [14] use on-surface feature consistency constraints, but the  
 437 reconstruction results are still not good enough. In con-  
 438 trast, our method achieves more comprehensive geometry  
 439 and finer details by employing volume rendering-based fea-  
 440 ture consistency constraints. This highlights the advantages  
 441 of our approach in geometric consistency. More visualiza-  
 442 tions are presented in the supplementary materials.

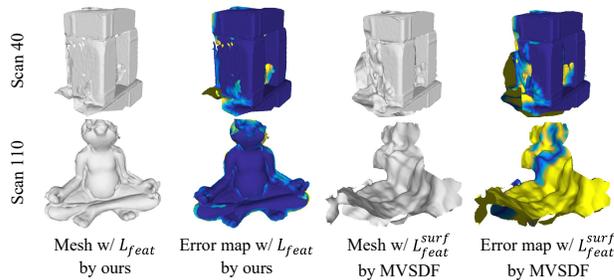


Figure 8. Reconstructed meshes and error maps on DTU dataset with different feature consistency losses.

### 4.5. Ablation Study

443 We evaluate the components of our method with 3 small-  
 444 overlapping views by an ablation study on the DTU [16]  
 445 dataset. To compare the depth loss  $L_{depth}^{mono}$  calibrated  
 446 by rendered depth in MonoSDF [46] with our depth loss  
 447  $L_{depth}$ , we replace  $L_{depth}$  with  $L_{depth}^{mono}$  to evaluate it in  
 448 our method. We also compare the volume rendering-based  
 449

Method	$L_{color}$	$L_{feat}$	$L_{depth}$	$L_{depth}^{mono}$	$L_{feat}^{L1}$	$L_{feat}^{L2}$	$L_{feat}^{surf}$	CD↓
Baseline								3.35
	✓							1.76
	✓	✓						1.47
	✓		✓					1.62
	✓	✓	✓					<b>1.11</b>
	✓			✓				1.59
	✓		✓		✓			2.36
	✓		✓			✓		1.81
	✓		✓				✓	2.93

Table 2. Ablation studies on DTU dataset with 3 small-overlapping images.

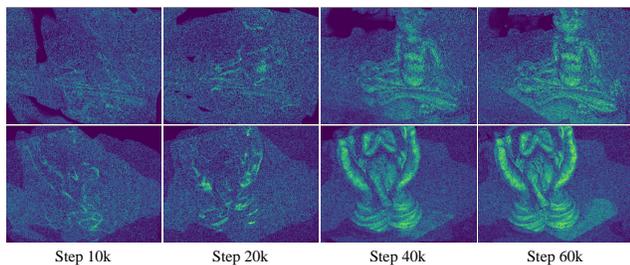


Figure 9. The variation of weighted feature similarity during training, brighter colors indicate higher feature similarity.

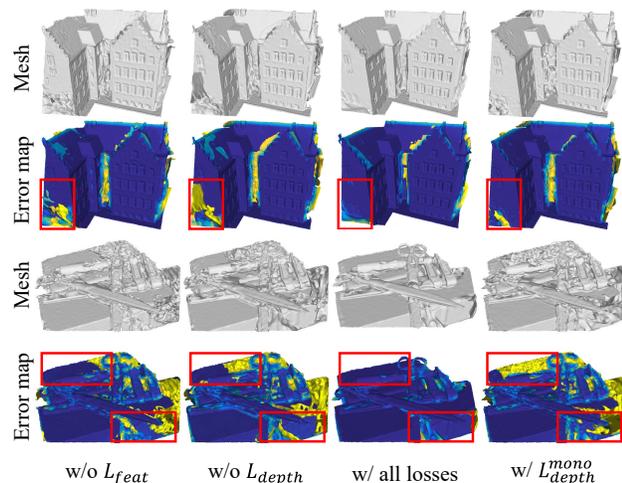


Figure 10. Visualization of reconstruction and error maps for scene scan24 and scan37 in DTU dataset with different losses. The differences of error maps are highlighted.

450 feature consistency loss calculated using L1 distance (de-  
451 noted as  $L_{feat}^{L1}$ ) and L2 distance (denoted as  $L_{feat}^{L2}$ ) with  
452 our method using feature similarity distance. We found that  
453 feature similarity distance is better than both L1 and L2 dis-  
454 tance, as shown in Table 2.

455 In addition, we replace our volume rendering-based fea-  
456 ture consistency loss  $L_{feat}$  with the on-surface feature con-  
457 sistency loss  $L_{feat}^{surf}$  used in MVSDf [50] to compare the ef-  
458 fects of two different loss functions. Figure 8 illustrates the

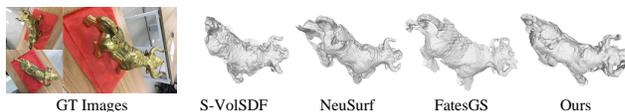


Figure 11. Failure case. For specular objects, the ambiguity in the color consistency constraint may lead to a rough surface.

reconstruction results and error maps on the DTU dataset  
when using different feature consistency losses, under the  
on-surface feature consistency loss  $L_{feat}^{surf}$ , the meshes show  
large artifacts.

Table 2 shows the average Chamfer Distance over all 11  
scenes on DTU dataset using different losses. The exper-  
imental results indicate that both feature consistency loss  
and uncertainty-guided depth constraint improve the surface  
reconstruction.

Figure 9 illustrates the variation of the weighted feature  
similarity map during the training process. Brighter colors  
indicate higher feature similarity, demonstrating that our  
volume rendering-based feature consistency loss can provide  
effective constraints.

Figure 10 shows the reconstructed meshes and error  
maps for scene scan24 and scan37 on the DTU [16] dataset  
when using different losses. It can be observed that the  
mesh deteriorates with out the volume rendering-based fea-  
ture consistency loss or the uncertainty-guided depth con-  
straint loss, and the reconstruction quality drops when using  
the depth loss  $L_{depth}^{mono}$  in MonoSDF [46].

## 5. Conclusions

We propose a novel method for learning implicit repre-  
sentations from sparse views with small overlaps. Our novelty  
lies in a novel volume rendering-based feature consistency  
loss and an uncertainty-guided depth constraint. Extensive  
experiments on the DTU [16] and BlendedMVS [43] data-  
sets show that our method surpasses existing state-of-  
the-art sparse-view reconstruction methods in terms of re-  
construction quality.

**Limitations.** Although our method shows significant  
improvement over other sparse view reconstruction meth-  
ods, there are still some limitations. Firstly, for specular  
objects, the ambiguity in the color consistency constraint  
may lead to a rough surface, as shown in Figure 11.  
Secondly, Following previous studies [13, 23, 39], the  
camera poses of sparse views are obtained from the training  
dataset. However, in some cases, it may not be possible  
to obtain accurate camera poses using SfM methods like  
COLMAP [31] due to the lack of texture in the images or  
excessive viewing angles. Additionally, the feature con-  
sistency constraint method requires a pre-trained network  
to extract image features. The accuracy of the features  
determines the performance of the feature consistency  
constraint.

504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559

References

[1] Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haomin Liu, Hujun Bao, and Guofeng Zhang. PGSR: Planar-based gaussian splatting for efficient and high-fidelity surface reconstruction. *IEEE Transactions on Visualization and Computer Graphics*, PP: 1–12, 2024. 3, 6

[2] Hanlin Chen, Chen Li, and Gim Hee Lee. NeuSG: Neural implicit surface reconstruction with 3D gaussian splatting guidance. *arXiv preprint arXiv:2312.00846*, 2023. 3

[3] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. MVSpLat: Efficient 3D gaussian splatting from sparse multi-view images. In *European Conference on Computer Vision*, pages 370–386. Springer, 2025. 3

[4] Pinxuan Dai, Jiamin Xu, Wenxiang Xie, Xinguo Liu, Huamin Wang, and Weiwei Xu. High-quality surface reconstruction using gaussian surfels. In *ACM SIGGRAPH 2024 Conference Proceedings*, 2024. 3

[5] François Darmon, Bénédicte Bascle, Jean-Clément Devaux, Pascal Monasse, and Mathieu Aubry. Improving neural implicit surfaces geometry with patch warping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6260–6269, 2022. 1, 2, 5, 6

[6] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022. 3

[7] Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. Geo-Neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 35:3403–3416, 2022. 2

[8] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Proceedings of Machine Learning and Systems 2020*, pages 3569–3579. 2020. 5

[9] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. 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*, pages 2495–2504, 2020. 2

[10] Antoine Guédon and Vincent Lepetit. SuGaR: Surface-aligned gaussian splatting for efficient 3D mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5354–5363, 2024. 3

[11] Liang Han, Junsheng Zhou, Yu-Shen Liu, and Zhizhong Han. Binocular-guided 3D gaussian splatting with view consistency for sparse view synthesis. In *Advances in Neural Information Processing Systems*, 2024. 3

[12] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2D gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 Conference Proceedings*, 2024. 3, 6

[13] Han Huang, Yulun Wu, Junsheng Zhou, Ge Gao, Ming Gu, and Yu-Shen Liu. NeuSurf: On-surface priors for neural surface reconstruction from sparse input views. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2312–2320, 2024. 1, 2, 4, 6, 7, 8

[14] Han Huang, Yulun Wu, Chao Deng, Ge Gao, Ming Gu, and Yu-Shen Liu. FatesGS: Fast and accurate sparse-view surface reconstruction using gaussian splatting with depth-feature consistency. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025. 1, 3, 4, 6, 7

[15] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting NeRF on a Diet: Semantically consistent few-shot view synthesis implementation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. 3

[16] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 406–413, 2014. 2, 5, 7, 8

[17] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning (ICML)*, pages 5156–5165, 2020. 2

[18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 3

[19] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. DNGaussian: Optimizing sparse-view 3D gaussian radiance fields with global-local depth normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20775–20785, 2024. 3

[20] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8456–8465, 2023. 2

[21] Yixun Liang, Hao He, and Yingcong Chen. ReTR: Modeling rendering via transformer for generalizable neural surface reconstruction. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 2, 6

[22] Xinqi Liu, Jituo Li, and Guodong Lu. Reconstructing complex shaped clothing from a single image with feature stable unsigned distance fields. *IEEE Transactions on Visualization and Computer Graphics*, 2024. 1

[23] Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. SparseNeuS: Fast generalizable neural surface reconstruction from sparse views. In *European Conference on Computer Vision*, pages 210–227. Springer, 2022. 1, 2, 6, 7, 8

[24] Xiaoyang Lyu, Yang-Tian Sun, Yi-Hua Huang, Xiuzhe Wu, Ziyi Yang, Yilun Chen, Jiangmiao Pang, and Xiaojuan Qi. 3DGSr: Implicit surface reconstruction with 3D gaussian splatting. *arXiv preprint arXiv:2404.00409*, 2024. 3

560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616

617 [25] Baorui Ma, Yu-Shen Liu, and Zhizhong Han. Learning  
618 signed distance functions from noisy 3D point clouds via  
619 noise to noise mapping. In *International Conference on Ma-*  
620 *chine Learning (ICML)*, 2023. 1

621 [26] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik,  
622 Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF:  
623 Representing scenes as neural radiance fields for view syn-  
624 thesis. *Communications of the ACM*, 65(1):99–106, 2021.  
625 2

626 [27] Youngju Na, Woo Jae Kim, Kyu Beom Han, Suhyeon Ha,  
627 and Sung-Eui Yoon. UFORecon: Generalizable sparse-view  
628 surface reconstruction from arbitrary and unfavorable sets.  
629 In *Proceedings of the IEEE/CVF Conference on Computer*  
630 *Vision and Pattern Recognition*, pages 5094–5104, 2024. 1,  
631 2, 6, 7

632 [28] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall,  
633 Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Reg-  
634 NeRF: Regularizing neural radiance fields for view synthesis  
635 from sparse inputs. In *Proceedings of the IEEE/CVF Con-*  
636 *ference on Computer Vision and Pattern Recognition*, pages  
637 5480–5490, 2022. 3

638 [29] Rui Peng, Xiaodong Gu, Luyang Tang, Shihe Shen, Fanqi  
639 Yu, and Ronggang Wang. GenS: Generalizable neural sur-  
640 face reconstruction from multi-view images. In *Advances*  
641 *in Neural Information Processing Systems*, pages 56932–  
642 56945, 2023. 1, 2, 6

643 [30] Yufan Ren, Tong Zhang, Marc Pollefeys, Sabine Süs-  
644 trunk, and Fangjinhua Wang. VolRecon: Volume rendering  
645 of signed ray distance functions for generalizable multi-view  
646 reconstruction. In *Proceedings of the IEEE/CVF Conference*  
647 *on Computer Vision and Pattern Recognition*, pages 16685–  
648 16695, 2023. 1, 2, 6

649 [31] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys,  
650 and Jan-Michael Frahm. Pixelwise view selection for un-  
651 structured multi-view stereo. In *European Conference on*  
652 *Computer Vision*, 2016. 2, 4, 5, 8

653 [32] Hui Tian, Chenyang Zhu, Yifei Shi, and Kai Xu. SuperUDF:  
654 Self-supervised udf estimation for surface reconstruction.  
655 *IEEE Transactions on Visualization and Computer Graph-*  
656 *ics*, 2023. 1

657 [33] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt,  
658 and Federico Tombari. SPARF: Neural radiance fields from  
659 sparse and noisy poses. In *Proceedings of the IEEE/CVF*  
660 *Conference on Computer Vision and Pattern Recognition*,  
661 pages 4190–4200, 2023. 3

662 [34] Guangcong Wang, Zhaoxi Chen, Chen Change Loy, and Zi-  
663 wei Liu. SparseNeRF: Distilling depth ranking for few-shot  
664 novel view synthesis. In *Proceedings of the IEEE/CVF In-*  
665 *ternational Conference on Computer Vision*, pages 9065–9076,  
666 2023. 3

667 [35] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku  
668 Komura, and Wenping Wang. NeuS: Learning neural im-  
669 plicit surfaces by volume rendering for multi-view recon-  
670 struction. *Advances in Neural Information Processing Sys-*  
671 *tems*, 2021. 1, 2, 3, 5, 6

672 [36] Yiqun Wang, Ivan Skorokhodov, and Peter Wonka.  
673 HF-NeuS: Improved surface reconstruction using high-  
frequency details. *Advances in Neural Information Process-*  
674 *ing Systems*, 35:1966–1978, 2022.

[37] Yiming Wang, Qin Han, Marc Habermann, Kostas Dani-  
675 ilidis, Christian Theobalt, and Lingjie Liu. NeuS2: Fast  
676 learning of neural implicit surfaces for multi-view recon-  
677 struction. In *Proceedings of the IEEE/CVF International*  
678 *Conference on Computer Vision*, pages 3295–3306, 2023. 1,  
679 2  
680 2  
681 2

[38] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-  
682 moncelli. Image quality assessment: from error visibility to  
683 structural similarity. *IEEE Transactions on Image Process-*  
684 *ing*, 13(4):600–612, 2004. 5

[39] Haoyu Wu, Alexandros Graikos, and Dimitris Samaras. S-  
685 VolSDF: Sparse multi-view stereo regularization of neural  
686 implicit surfaces. *International Conference on Computer Vi-*  
687 *sion*, 2023. 1, 2, 6, 7, 8

[40] Tong Wu, Jiaqi Wang, Xingang Pan, Xudong Xu, Christian  
688 Theobalt, Ziwei Liu, and Dahua Lin. Voxurf: Voxel-based  
689 efficient and accurate neural surface reconstruction. *Inter-*  
690 *national Conference on Learning Representations*, 2022. 1,  
691 2  
692 2

[41] Luoyuan Xu, Tao Guan, Yuesong Wang, Wenkai Liu, Zhao-  
693 jie Zeng, Junle Wang, and Wei Yang. C2F2NeUS: Cascade  
694 cost frustum fusion for high fidelity and generalizable neu-  
695 ral surface reconstruction. In *Proceedings of the IEEE/CVF*  
696 *International Conference on Computer Vision*, pages 18291–  
697 18301, 2023. 2

[42] Jiawei Yang, Marco Pavone, and Yue Wang. FreeNeRF: Im-  
698 proving few-shot neural rendering with free frequency reg-  
699 ularization. In *Proceedings of the IEEE/CVF Conference*  
700 *on Computer Vision and Pattern Recognition*, pages 8254–  
701 8263, 2023. 3

[43] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan  
702 Ren, Lei Zhou, Tian Fang, and Long Quan. Blended-  
703 MVS: A large-scale dataset for generalized multi-view stereo  
704 networks. In *Proceedings of the IEEE/CVF Conference*  
705 *on Computer Vision and Pattern Recognition*, pages 1790–  
706 1799, 2020. 2, 5, 6, 7, 8

[44] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Vol-  
707 ume rendering of neural implicit surfaces. *Advances in Neu-*  
708 *ral Information Processing Systems*, 34:4805–4815, 2021. 1,  
709 2, 6  
710 2

[45] Mae Younes, Amine Ouasfi, and Adnane Boukhayma. Spar-  
711 seCraft: Few-shot neural reconstruction through stereopsis  
712 guided geometric linearization. In *European Conference on*  
713 *Computer Vision*, pages 37–56. Springer, 2024. 1, 6, 7

[46] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sat-  
714 tler, and Andreas Geiger. MonoSDF: Exploring monocu-  
715 lar geometric cues for neural implicit surface reconstruc-  
716 tion. *Advances in Neural Information Processing Systems*,  
717 35:25018–25032, 2022. 1, 2, 4, 6, 7, 8

[47] Zehao Yu, Torsten Sattler, and Andreas Geiger. Gaussian  
718 Opacity Fields: Efficient adaptive surface reconstruction in  
719 unbounded scenes. *ACM Transactions on Graphics*, 43(6):  
720 1–13, 2024. 3

[48] Yu-Jie Yuan, Yu-Kun Lai, Yi-Hua Huang, Leif Kobbelt, and  
721 Lin Gao. Neural radiance fields from sparse rgb-d images for  
722 723 724 725 726 727 728 729 730

- 731 high-quality view synthesis. *IEEE Transactions on Pattern*  
732 *Analysis and Machine Intelligence*, 45(7):8713–8728, 2022.  
733 3
- 734 [49] Jingyang Zhang, Yao Yao, Shiwei Li, Zixin Luo, and Tian  
735 Fang. Visibility-aware multi-view stereo network. *The*  
736 *British Machine Vision Conference*, 2020. 2, 3, 5, 6
- 737 [50] Jingyang Zhang, Yao Yao, and Long Quan. Learning signed  
738 distance field for multi-view surface reconstruction. In  
739 *Proceedings of the IEEE/CVF International Conference on*  
740 *Computer Vision*, pages 6525–6534, 2021. 4, 6, 8
- 741 [51] Jiawei Zhang, Jiahe Li, Xiaohan Yu, Lei Huang, Lin Gu,  
742 Jin Zheng, and Xiao Bai. CoR-GS: sparse-view 3d gaussian  
743 splatting via co-regularization. In *European Conference on*  
744 *Computer Vision*, pages 335–352. Springer, 2024. 3
- 745 [52] Wenyuan Zhang, Ruofan Xing, Yunfan Zeng, Yu-Shen Liu,  
746 Kanle Shi, and Zhizhong Han. Fast learning radiance fields  
747 by shooting much fewer rays. *IEEE Transactions on Image*  
748 *Processing*, 2023. 1
- 749 [53] Wenyuan Zhang, Yu-Shen Liu, and Zhizhong Han. Neu-  
750 ral signed distance function inference through splatting 3d  
751 gaussians pulled on zero-level set. In *Advances in Neural*  
752 *Information Processing Systems*, 2024. 3
- 753 [54] Junsheng Zhou, Baorui Ma, Shujuan Li, Yu-Shen Liu, Yi  
754 Fang, and Zhizhong Han. CAP-UDF: Learning unsigned  
755 distance functions progressively from raw point clouds with  
756 consistency-aware field optimization. *IEEE Transactions on*  
757 *Pattern Analysis and Machine Intelligence*, 2024. 1, 2
- 758 [55] Junsheng Zhou, Baorui Ma, and Liu Yu-Shen. Fast learn-  
759 ing of signed distance functions from noisy point clouds via  
760 noise to noise mapping. *IEEE Transactions on Pattern Anal-*  
761 *ysis and Machine Intelligence*, 2024. 1
- 762 [56] Tiansong Zhou, Jing Huang, Tao Yu, Ruizhi Shao, and Kun  
763 Li. HDhuman: High-quality human novel-view rendering  
764 from sparse views. *IEEE Transactions on Visualization and*  
765 *Computer Graphics*, 2023. 1
- 766 [57] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang.  
767 FSGS: Real-time few-shot view synthesis using gaussian  
768 splatting. In *European Conference on Computer Vision*,  
769 pages 145–163. Springer, 2025. 3