

GeoTransfer: Generalizable Few-Shot Multi-View Reconstruction via Transfer Learning

Shubhendu Jena, Franck Multon, and Adnane Boukhayma

Inria, Univ. Rennes, CNRS, IRISA, M2S, France

Abstract. This paper presents a novel approach for sparse 3D reconstruction by leveraging the expressive power of Neural Radiance Fields (NeRFs) [35] and fast transfer of their features to learn accurate occupancy fields. Existing 3D reconstruction methods from sparse inputs still struggle with capturing intricate geometric details and can suffer from limitations in handling occluded regions. On the other hand, NeRFs [35] excel in modeling complex scenes but do not offer means to extract meaningful geometry. Our proposed method offers the best of both worlds by transferring the information encoded in NeRF [35] features to derive an accurate occupancy field representation. We utilize a pre-trained, generalizable state-of-the-art NeRF network [18] to capture detailed scene radiance information, and rapidly transfer this knowledge to train a generalizable implicit occupancy network. This process helps in leveraging the knowledge of the scene geometry encoded in the generalizable NeRF [18] prior and refining it to learn occupancy fields, facilitating a more precise generalizable representation of 3D space. The transfer learning approach leads to a dramatic reduction in training time, by orders of magnitude (*i.e.* from several days to 3.5 hrs), obviating the need to train generalizable sparse surface reconstruction methods from scratch. Additionally, we introduce a novel loss on volumetric rendering weights that helps in the learning of accurate occupancy fields, along with a normal loss that helps in global smoothing of the occupancy fields. We evaluate our approach on the DTU dataset [1] and demonstrate **state-of-the-art performance** in terms of reconstruction accuracy, especially in challenging scenarios with sparse input data and occluded regions. We furthermore demonstrate the generalization capabilities of our method by showing qualitative results on the Blended MVS [59] dataset without any re-training. Project page : <https://shubhendu-jena.github.io/geotransfer/>

Keywords: 3D Reconstruction · Volume rendering · Sparse views

1 Introduction

Creating three-dimensional structures from a set of images is a fundamental task in the realm of computer vision, finding broad applications in fields like robotics, augmented reality, and virtual reality. The first seminal deep learning approaches tackling this problem used Multi-view stereo techniques (MVS), as demonstrated by MVSNet [58] and its successors [10, 14, 50, 57]. These methods build 3D cost volumes based on the camera frustum, deviating from traditional euclidean space, to achieve accurate depth-map estimation. However, they often

necessitate subsequent steps, such as depth-map filtering, fusion, and mesh reconstruction, and exhibit susceptibility to noise, texture-less regions, and gaps. Unlike seminal work relying on explicit representations (*e.g.* meshes [16, 19, 51] and point clouds [2, 12, 21]), neural implicit reconstruction methods [8, 37, 38, 52, 60] constitute another popular class of strategies to address this challenge, creating precise and realistic geometry from multi-view images through the use of volume rendering and neural implicit representations based on the Sign Distance Function (SDF) [44] and its variations. However, despite their effectiveness, these approaches come with inherent limitations, including a lack of cross-scene generalization capabilities and the need for substantial computational resources to train them from scratch for each scene. Moreover, these techniques heavily depend on a large number of input views. However, due to many constrained scenarios (*e.g.* out-of-the-studio, low budget, *etc.*) and in the interest of wider applicability, there is active interest in seeking solutions that can deliver under minimal input.

To solve these issues, recent investigations, in the context of novel-view synthesis [18, 36, 56] and 3D reconstruction [26, 29, 47] have sought solutions by relying 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. This approach has proven effective in achieving remarkable cross-scene generalization capabilities, even with sparse views as input. Most relevant to our approach is GeoNeRF [18] which constructs a cost volume to enable geometry-aware scene reasoning, followed by attention-based view aggregation and volumetric rendering to learn the radiance field of a scene. In this paper, differently from concurrent works, we explore the idea of using this pre-existing generalizable state-of-the-art NeRF to obtain scene reconstructions through transfer learning. In this context, we show that under the assumption of our scene being composed of solid, non-transparent objects, it is possible to rapidly transform the generalizable sampling-dependent opacity obtained from the density field of GeoNeRF [18] to a generalizable sampling-independent occupancy. This strategy also leads to a drastic reduction of training time from the order of several days to a couple of hours, which removes the need to train generalizable sparse reconstruction methods [26, 29, 47] from scratch. We introduce a novel volumetric rendering weight loss, in addition to a surface normal based smoothing loss to further refine our occupancy field, leading to the current state-of-the-art results in sparse 3D reconstruction on the DTU dataset [1] without requiring any test-time optimization, and outperforming in this process the current state-of-the-art generalizable SDF based 3D reconstruction networks [26, 29, 47]. In summary, our contribution can be summarized as:

- We explore a novel strategy of fast adaptation of an existing state-of-the-art generalizable NeRF method to obtain a generalizable occupancy network by transferring and fine-tuning its features. This yields the **state-of-the-art performance on DTU [1]** reconstruction from sparse views 1.
- The decrease in training duration from multiple days to just a few hours, all the while achieving state-of-the-art performance, removes the necessity to

train computationally intensive sparse generalizable surface reconstruction techniques from scratch.

- Among the losses we use for our transfer learning framework, we propose a novel volumetric rendering weight loss to impose the properties followed by an ideal occupancy field in a volumetric rendering framework which leads to the learning of a more accurate occupancy functions (5).

2 Related Work

There is a substantial body of work on the subject of 3D reconstruction, and we review here work we deemed most relevant to the context of our contribution.

Neural Surface Reconstruction. In the realm of neural surface reconstruction, the utilization of neural implicit representations enables the depiction of 3D geometries as continuous functions that can be computed at arbitrary spatial locations. Due to their capability to represent complex and detailed shapes in a compact and efficient manner, these representations demonstrate significant potential in tasks such as 3D reconstruction [8, 17, 20, 37, 38, 52, 60–62, 64], shape representation [3, 13, 33, 44], and novel view synthesis [27, 35, 49]. The emergence of NeRF [35] has instigated a significant shift in the paradigm towards employing similar techniques for these tasks. IDR [61] utilizes surface rendering to acquire geometry from multi-view images but necessitates additional object masks. Unisurf [38], which is most relevant to our work, models the local opacity of a NeRF [35] with an occupancy network, which allows them to train on multiple datasets including ones involving forward-facing scenes such as LLFF [34]. Differently, several methods have attempted to rewrite the density function in NeRF [35] using Signed Distance Function (SDF) and its variants, successfully achieving plausible geometry. Significantly, NeuS [52] formulates an unbiased and occlusion aware volumetric weight function equation by employing logistic sigmoid functions. Conversely, Volsdf [60] incorporates a signed distance function into the density formulation and introduces a sampling strategy that meets a determined error bound on the transparency function. HF-NeuS [54] improves upon NeuS [52] by modeling transparency as a transformation of the estimated signed distance field and proposes to decompose the signed distance function into a base function and a displacement function with a coarse-to-fine strategy to gradually increase the high-frequency details. These methods offer a robust approach for multi-view 3D reconstruction from 2D images. However, these methods require prolonged optimization for training each scene independently and also require a substantial number of dense images, making it challenging to generalize to unknown scenes and limiting deployment.

Generalizable NeRFs. Some recent methods [9, 15, 22, 25, 36, 55, 62] synthesize novel views on a single scene with sparse views, albeit facing difficulties in understanding the underlying geometry of the scene, which they attempt to solve using several geometry-based regularization strategies [4, 13, 40, 43]. To attempt to solve this problem, certain methods [6, 7, 18, 24, 28, 53, 63] generate novel views

in unknown scenarios through a generalization approach, constructing neural radiance fields on sparse views. These methods can infer on unknown scenarios without fine-tuning after training on multiple known scenarios, which involves incorporating priors derived from a larger model trained on multi-view image datasets. Among such methods, PixelNeRF [63] conditions itself on features extracted by a CNN. MVSNeRF [6] constructs a neural volume from the cost volume obtained by warping image features and conditions itself on this neural volume. IBRNet [53] aggregates features from nearby views to infer geometry and adopts an image-based rendering approach. NeuRay [28] utilizes neural networks to model and handle occlusions, thereby improving the quality and accuracy of image-based rendering. GeoNeRF [18], another recent method and one that we build on, employs a cascaded cost volume and an attention-based technique to aggregate information from different views. However, deriving the scene geometry from the volume density of these NeRF [35] based methods involves meticulous tuning of the density threshold, leading to artifacts due to the inherent ambiguity in the density field, as previously pointed out in Unisurf [38].

Generalizable Neural Surface Reconstruction. Point Cloud input models [5, 39, 41, 42, 46] are typically trained with SDF ground truth supervision. To obtain 3D reconstructions from sparse images with an accurate level set, current methods marry generalizable NeRFs with the Signed Distance Function (SDF) based transformation functions that model the volume density, thereby enabling volume rendering. Among these methods, both SparseNeuS [29] and VolRecon [47] achieve generalizable neural surface reconstruction by utilizing information from source images as priors. SparseNeuS [29] does this by encoding geometric information using a regular euclidean volume, while VolRecon [47] introduces multi-view image features through the view transformer to advance this scheme. The most recent work of ReTR [26] utilizes a hybrid extractor to obtain a multi-level euclidean volume and then employs a reconstruction transformer to enhance performance. Contrarily to these methods, we employ the Unisurf [38] based volumetric rendering framework which reconstructs surfaces by predicting occupancy, enabling the seamless transfer of the opacity information from a pretrained GeoNeRF [18] to learn a refined occupancy field.

3 Method

Our goal is to obtain a feed-forward generalizable 3D reconstruction network from a sparse array of images. Namely, this network (*e.g.* f_o) should be able to deliver an implicit shape representation, *i.e.* a binary shape occupancy field for instance, from observed images $\{I_i\}_{i=1}^N$ and their calibrations $\{\pi_i\}_{i=1}^N$ of a test scene (unseen at training), without requiring any optimization on this new scene. The inferred shape $\hat{\mathcal{S}}$ can be obtained as the level set of the learned occupancy f_o :

$$\hat{\mathcal{S}} = \{\mathbf{x} \in \mathbb{R}^3 \mid f_o(\mathbf{x}) = 0.5\}. \quad (1)$$

Practically, an explicit triangle mesh for $\hat{\mathcal{S}}$ can be obtained through the Marching Cubes algorithm [30] while querying the neural network f_o . We propose to adapt

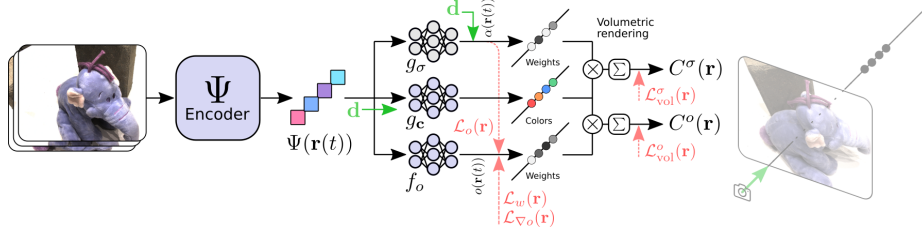


Fig. 1: Overview of our transfer learning: Our final model (in light blue) is comprised of tuned image encoder Ψ and implicit occupancy and color decoders f_o and g_c . The encoder Ψ , the color decoder g_c and density decoder g_σ are initialized as a pretrained generalizable NeRF. Red dashed lines symbolize our tuning losses. We apply multiple regularizations on our occupancy f_o , while tuning the network with both the density and occupancy guided volumetric renderings.

an existing generalizable NeRF model for this purpose. We chose the model denoted GeoNeRF [18] without loss of generality as it is one of the best performing models in generalizable novel view synthesis.

Generalizable NeRF model Let \mathbf{r} be a ray, *i.e.* $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$, 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 by integrating along the ray:

$$C(\mathbf{r}) = \int T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt \quad (2)$$

$$= \int w(t)\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt. \quad (3)$$

In the equation above, $\sigma(\mathbf{r}(t))$ 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:

$$T(t) = \exp\left(-\int_0^t \sigma(\mathbf{r}(s))ds\right). \quad (4)$$

Furthermore, $w(t)$ is commonly referred to as the volumetric rendering weighting function.

In practice, the integral in Equation 3 is approximated using discrete samples $\{t_i\}$ with the quadrature rule [32]. This gives rise to an equation resembling α -compositing, where α represents opacity, which controls the amount of radiance that is absorbed or transmitted at each point in the scene. This leads to :

$$C_\sigma(\mathbf{r}) = \sum_i T_i \left(1 - e^{-\sigma_i(t_{i+1}-t_i)}\right) \mathbf{c}(\mathbf{r}(t_i), \mathbf{d}) \quad (5)$$

$$= \sum_i \prod_{j=0}^{i-1} (1 - \alpha_j) \alpha_i \mathbf{c}(\mathbf{r}(t_i), \mathbf{d}) \quad (6)$$

where $\alpha_i = 1 - \exp(-\sigma_i(t_{i+1} - t_i))$.

Generalizable NeRFs are typically comprised of an encoder network (*e.g.* Ψ) producing spatially local features. These features are mapped subsequently by a density network (MLP) (*e.g.* g_σ) and a viewing direction dependent color network (MLP) (*e.g.* g_c). (*cf.* Figure 1). Inference is achieved via volumetric rendering as shown above. The colors and densities necessary for this rendering are thus modelled as follows:

$$\mathbf{c}(\mathbf{r}(t), \mathbf{d}) = g_\sigma(\Psi(\mathbf{r}(t), \{I_i\}), \mathbf{d}) \quad (7)$$

$$\sigma(\mathbf{r}(t)) = g_c(\Psi(\mathbf{r}(t), \{I_i\})). \quad (8)$$

Adapting a Generalizable NeRF model A geometry can be extracted from NeRF [35] models at test time by thresholding the density function, which is view and sampling independent. However, it is not clear how the threshold can be selected, and the obtained geometries tend to be noisy and inaccurate, with high Chamfer errors with respect to the ground truth.

Although the opacity is the closest measure in a NeRF [35] to an occupancy, the opacity as defined in the volumetric rendering framework is view and ray-sampling dependent. Hence it is not clear how a solely spatially dependent geometry can be extracted from it.

Yet, in order for generalizable NeRFs to be able to reason correctly about 3D for accurate novel view synthesis, we hypothesize that they must encompass a good high level representation of geometry, that could be nudged towards an accurate and smooth shape representation per se. Based on this, we propose to tune a Generalizable NeRF (Namely GeoNeRF [18]) into a generalizable occupancy model, that offers view and ray-sampling independent geometry at convergence.

We define a new Sigmoid activated implicit decoder f_o that will represent the occupancy field in feature space:

$$o(\mathbf{r}(t)) = f_o(\Psi(\mathbf{r}(t), \{I_i\})) \quad (9)$$

We propose to tune the encoder Ψ to adapt the feature space to this new prediction task. This can be seen as form of transfer learning. As this tuning is achieved through volumetric rendering, we use the baseline model color network ($\mathbf{c} := g_c$) to perform volumetric rendering with f_o :

$$C_o(\mathbf{r}) = \sum_i \prod_{j=0}^{i-1} (1 - o_j) o_i \mathbf{c}(\mathbf{r}(t_i), \mathbf{d}). \quad (10)$$

Color network g_c has been pretrained using the density based weight functions in volumetric rendering. Hence, it is important to tune it as-well to this new occupancy based rendering. Furthermore, as we also want our initial Generalizable NeRF not to deviate substantially from its original weights and hence lose its original knowledge, we train by backpropagating both the volumetric rendering loss based on occupancy, and the original volumetric rendering loss

based on density together:

$$\mathcal{L}_{vol}^o(\mathbf{r}) = \|C_o(\mathbf{r}) - C_{GT}(\mathbf{r})\|_2 \quad (11)$$

$$\mathcal{L}_{vol}^\sigma(\mathbf{r}) = \|C_\sigma(\mathbf{r}) - C_{GT}(\mathbf{r})\|_2. \quad (12)$$

We also retain the depth supervision used by the geometry reasoner in GeoNeRF [18], with ground truth depths if available, and with their self-supervised depth loss otherwise. Based on the knowledge that assuming solid objects, α becomes a discrete occupancy indicator variable $o \in \{0, 1\}$ which either takes $o = 0$ in free space and $o = 1$ in occupied space, we bootstrap our occupancy with α predictions from the density branch, as a form of warm up or initialization for a few iterations at the beginning of the training:

$$\mathcal{L}_o(\mathbf{r}) = \sum_i \|o_i - \alpha_i\|_2. \quad (13)$$

A common problem with learning geometry through volumetric rendering is that the weight function does not peak at the surface [38, 52]. This can be observed in Figure 5, where neither our baseline nor the generalizable NeRF model we build on satisfy this constraint. Hence, we propose a novel loss to remedy this limitation. First, for the current ray \mathbf{r} , we perform ray tracing, *i.e.* finding sample pair t_k and t_{k+1} where the occupancy flips from "empty" ($f_o(t_k) < 0.5$) to "occupied" ($f_o(t_{k+1}) \geq 0.5$) for the first time along the ray. Then, we perform secant method based root finding [37, 38] between these samples to find the root t^* corresponding to the surface-ray intersection. Ideally, we want our weight function to reach a sharp 1-peak at this location. Hence, we supervise the weight with a Gaussian centered at root t^* , and whose standard deviation we dynamically reduce during training, following the scheduling detailed in Section 4:

$$\mathcal{L}_w(\mathbf{r}) = \sum_t \left\| \prod_{j=0}^{i-1} (1 - o_j) o_i - e^{-((t_i - t^*)/a)^2} \right\|_2. \quad (14)$$

To reduce the noise in our reconstructions, we follow [38] and implement a smoothing over the normalized spatial gradients of our occupancy function at the surface, *i.e.* using root locations t^* :

$$\mathcal{L}_{\nabla o}(\mathbf{r}) = \left\| \frac{\nabla o(\mathbf{r}(t^*))}{\|\nabla o(\mathbf{r}(t^*))\|_2} - \frac{\nabla o(\mathbf{r}(t^*) + \epsilon)}{\|\nabla o(\mathbf{r}(t^*) + \epsilon)\|_2} \right\|_2, \quad (15)$$

where $\epsilon \in \mathbb{R}^3$ is a small random perturbation, and gradients can be computed efficiently through automatic differentiation [45].

Finally, our full training is done following the combined objective averaged over batches of rays, and we train on the same multi-view data as our original Generalizable NeRF model was initially trained on:

$$\mathcal{L} = \sum_{\mathbf{r}} \mathcal{L}_{vol}^o(\mathbf{r}) + \mathcal{L}_{vol}^\sigma(\mathbf{r}) + \lambda \mathcal{L}_{\nabla o}(\mathbf{r}) + \mu \mathcal{L}_w(\mathbf{r}) + \nu \mathcal{L}_o(\mathbf{r}). \quad (16)$$

These losses are analysed separately in section 4, where we show their respective numerical and qualitative contribution to our overall performance.

Scan	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122	Mean
COLMAP [48]	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
MVSNet [58]	1.05	2.52	1.71	1.04	1.45	1.52	0.88	1.29	1.38	1.05	0.91	<u>0.66</u>	0.61	1.08	1.16	1.22
IDR [61]	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 [60]	4.03	4.21	6.12	<u>0.91</u>	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 [38]	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 [52]	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
IBRNet-ft [53]	1.67	2.97	2.26	1.56	2.52	2.30	1.50	2.05	2.02	1.73	1.66	1.63	1.17	1.84	1.61	1.90
SparseNeuS-ft [29]	1.29	2.27	1.57	0.88	1.61	1.86	1.06	<u>1.27</u>	1.42	1.07	0.99	0.87	0.54	1.15	1.18	1.27
PixelNeRF [63]	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
IBRNet [53]	2.29	3.70	2.66	1.83	3.02	2.83	1.77	2.28	2.73	1.96	1.87	2.13	1.58	2.05	2.09	2.32
MVSNeRF [6]	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
GeoNeRF [18]	3.40	4.37	3.99	2.94	5.08	4.50	3.42	4.68	4.54	4.05	3.47	3.23	3.34	3.57	3.63	3.88
SparseNeuS [29]	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 [47]	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 [26]	1.05	2.31	1.44	0.98	1.18	1.52	0.88	1.35	<u>1.3</u>	0.87	1.07	0.77	0.59	1.05	1.12	1.17
Ours	1.01	<u>2.24</u>	1.52	0.88	1.37	1.82	<u>0.85</u>	1.39	1.25	<u>1.0</u>	0.77	<u>0.63</u>	<u>0.57</u>	<u>0.96</u>	1.0	<u>1.15</u>
<i>Ours*</i>	<u>0.95</u>	2.23	<u>1.45</u>	0.94	<u>1.26</u>	<u>1.67</u>	0.81	1.21	1.34	1.02	<u>0.84</u>	0.6	0.58	0.94	<u>1.02</u>	1.12

Table 1: Quantitative comparison on the DTU dataset [1]. Best and second best methods are **emboldened** and underlined respectively. *Ours** refers to the model trained using additional datasets.

4 Experiments

In this section, we showcase the efficacy and merits of our proposed approach. Firstly, we offer an intricate overview of our experimental configurations, encompassing implementation specifics, datasets, and baseline methods. Secondly, we present both quantitative and qualitative comparisons on two extensively utilized datasets, namely DTU [1] and BlendedMVS [59]. Finally, we perform thorough ablation studies to scrutinize the impact of distinct components in our proposed methodology.

Datasets In line with prior research [26, 29, 47], we employ the DTU dataset [1] for the training phase. Since we perform transfer learning from GeoNeRF [1], which is a NeRF-based framework, to learn generalizable occupancy fields, we are also able to leverage training on non-object centric datasets such as the real forward-facing datasets from LLFF [34] and IBRNet [53] and hence, we also include an additional model to show the effect of training on more data. In both cases (with and without additional training datasets), our backbone GeoNeRF [1] was trained on the same data as our full models to maintain fairness in evaluation performance comparisons. The DTU dataset [1] is characterized by indoor multi-view stereo data, featuring ground truth point clouds from 124 distinct scenes and under 7 different lighting conditions. Throughout our experiments, we utilize the same set of 15 scenes as [26, 29, 47] for testing purposes, reserving the remaining scenes for training. Concerning the BlendedMVS dataset [59], we opt for 7 scenes in accordance with SparseNeuS [29]. For each scene, we use the same 3 sparse input views following SparseNeuS [29]. To ensure impartial evaluation, we use the foreground masks from IDR [61] to evaluate how well our approach performs on the test set, consistent with prior research [26, 29, 47].

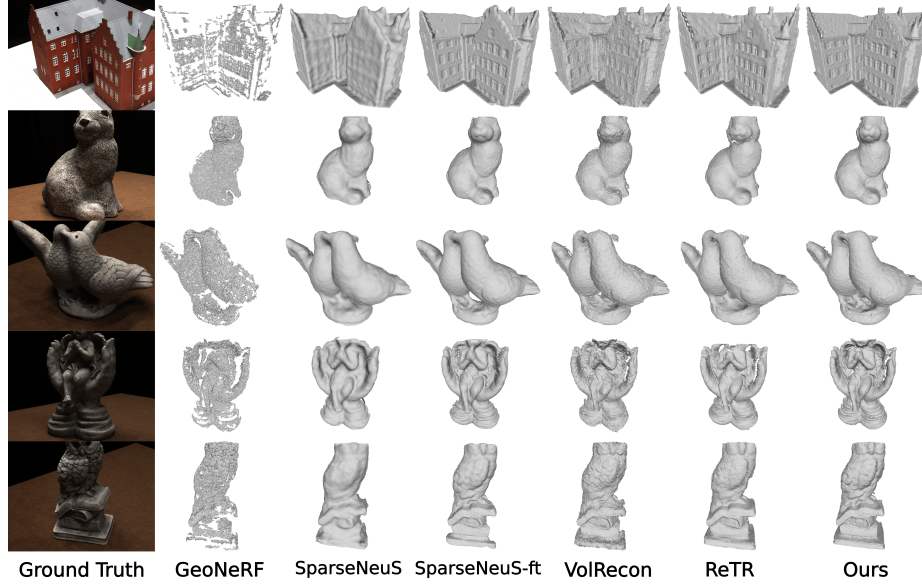


Fig. 2: Qualitative comparison of reconstructions from 3 input views in dataset DTU.

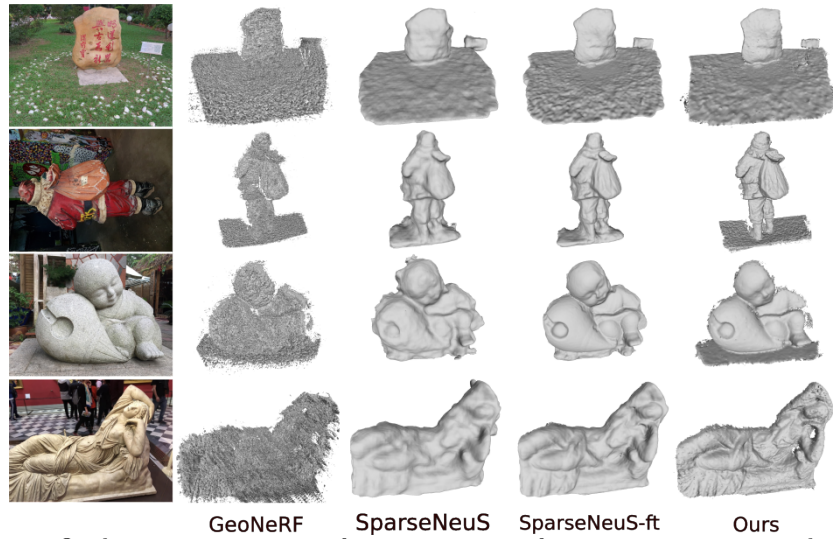


Fig. 3: Qualitative comparison of reconstructions from 3 input views in dataset BMVS. Note that we reconstruct detailed surfaces with our method without any fine-tuning.

Additionally, to assess the generalization capability of our proposed framework, we qualitatively compare our method on the BlendedMVS dataset [59] without any fine-tuning. For our novel-view synthesis experiments, we follow the same split of testing images within a scene as in GeoNeRF [1] and the testing scenes are identical to our 3D reconstruction experiments.

Baselines In order to showcase the efficacy of our method, we conducted comparisons with a) SparseNeus [29], VolRecon [47] and ReTR [26], the leading generalizable neural surface reconstruction approaches and their finetuned (ft) versions on sparse images; b) Generalizable neural rendering methods [6, 18, 53, 63]; c) Neural implicit reconstruction methods [38, 52, 60, 61] that necessitate scene-specific training from the beginning; and finally, d) Well-known multi-view stereo (MVS) [48, 58] methods.

Implementation details Our model is implemented using PyTorch [45] and PyTorch Lightning [11]. During the training phase, we utilize an image resolution of 800×600 , with N (the number of source images) set to 3. Training extends over 5400 steps using the Adam optimizer [23] on a single RTX A6000 GPU, with an initial learning rate of 5×10^{-4} . We apply the occupancy distillation loss for the first 100 warm-up steps. Thereafter, the weight of the distillation loss is reduced to 0. A cosine learning rate scheduler [31] is applied to the optimizer. The ray number sampled per batch and the batch size are configured to 128 and 1, respectively. Employing a hierarchical sampling strategy, we uniformly sample N_{coarse} points on the ray during both training and testing. Subsequently, importance sampling is applied to sample an additional N_{fine} points on top of the coarse probability estimation. In our experiments, we set N_{coarse} to 96 and N_{fine} to 32. During testing, the image resolution is kept the same at 800×600 . We also address the scheduling strategy followed for \mathcal{L}_w , described in 14. Specifically, we follow :

$$a = \max(a_{max}e^{-k\beta}, a_{min}). \quad (17)$$

where k denotes the global iteration number. a_{max} , a_{min} and β are hyperparameters, which we set to 1, 0.04 and 0.001 respectively. Finally, the weights associated with the losses are set as $\lambda = 0.1$, $\mu = 0.2$ and $\nu = 0.2$. We use the Marching Cubes algorithm [30] with a grid resolution of 400 for extracting each surface mesh by thresholding the occupancy field at 0.5.

4.1 Sparse View Reconstruction on DTU

We conduct surface reconstruction with sparse views (only 3 views) on the DTU dataset [1] and assess the predicted surface by comparing it to the ground-truth point clouds using the chamfer distance metric. To facilitate a fair comparison, we followed the evaluation process employed in [26, 29, 47] and adhered to the same testing split as described in them. As indicated in Table 1, our method (only DTU trained) surpasses SparseNeus [29] and VolRecon [47] by a considerable margin, *i.e.* by 30% and 17% respectively. When we use additional datasets for training, the gap further increases to 32% and 19% respectively. Furthermore,

our approach exhibits superior performance compared to well-known multi-view stereo (MVS) methods like Colmap [48] and MVSNet [58]. Our approach also demonstrates superior performance in comparison to ReTR [26], which is the latest state-of-the-art generalizable neural implicit reconstruction method. Additionally, we present qualitative results of sparse view reconstruction in Fig. 2, showcasing that our reconstructed geometry features more expressive and detailed surfaces which represent the ground truth surfaces more accurately when compared to the current state-of-the-art methods. Another DTU generalization experiment using MVSNeRF [6] as the backbone is included in the supplementary section to illustrate that the method presented can be extended to other generalizable NeRF baselines.

4.2 Generalization on BlendedMVS

To demonstrate the generalization prowess of our proposed approach, we perform additional evaluations on the BlendedMVS dataset [59] without resorting to any fine-tuning. The high-fidelity reconstructions of large-scale scenes and small objects across diverse domains, as illustrated in Fig. 3, affirms the efficacy of our method in terms of its generalization capabilities. Our method is able to obtain more detailed surfaces in comparison to SparseNeuS [29], even after it is fine-tuned on the sparse testing set.

4.3 Novel-view synthesis performance on DTU

Since GeoNeRF [18] is among the state-of-the-art methods for sparse novel-view synthesis, it is reasonable to assume that we also inherit its novel-view synthesis performance since we train our method by transferring its features. To validate this, we evaluate novel-view synthesis results on the DTU [1] dataset. We compare our method against GeoNeRF [18], VolRecon [47] and ReTR [26] on the mean PSNR, SSIM and LPIPS metrics over all scenes. Since we use 3 input source images to build our cost feature volume, for fairness in comparison, we also evaluate the novel-view synthesis results for VolRecon [47] and ReTR [26] with 3 input source images to build the cost feature volume, instead of the 4 input source images used during their respective trainings. Our results can be summarized as below :

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
VolRecon [47]	19.61	0.81	0.23
ReTR [26]	19.53	0.79	0.24
GeoNeRF [18]	23.48	0.93	0.087
Ours	24.08	0.93	0.093

Table 2: Novel-view performance on DTU [1].

As indicated in Table 2, we are very close in performance on novel-view synthesis metrics to GeoNeRF [18]. Note that we use 3 input views, the same

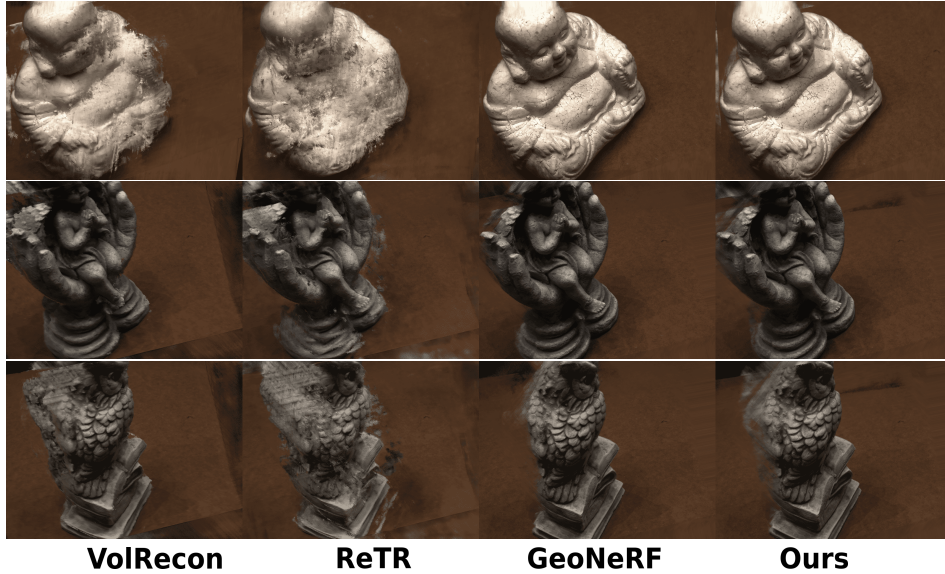


Fig. 4: Novel-View synthesis qualitative evaluation on DTU [1] using 3 source images. We notice that generalizable reconstruction models (ReTR [26], VolRecon [47]) struggle to perform extreme novel-view extrapolation.

setting on which we conduct our novel-view experiments on. The original GeoNeRF [18], however, was trained with 6 input views. The purpose behind this experiment was simply to make sure that our training paradigm does not sacrifice GeoNeRF [18]’s excellent novel-view synthesis performance. Overall, our method offers state-of-the-art reconstruction performance without foregoing the novel-view synthesis performance of one of the state-of-the-art novel-view synthesis methods we transfer features from. This is because our learnt occupancy fields by design and by penalization with the help of occupancy loss \mathcal{L}_o induce volumetric weights associated with the opacity fields of GeoNeRF [18]. We also notice that we far surpass VolRecon [47] and ReTR [26] in terms of novel-view synthesis performance. Hence, we offer a dual advantage, leading to our method offering **state-of-the-art** reconstruction results along with superior novel-view synthesis results. We also show some qualitative comparisons in Figure 4 and several other qualitative comparisons on several DTU [1] scenes in the supplementary section.

4.4 Ablation studies

In this section, we conduct an ablative analysis to justify the choice of our final architecture and choice of losses. We ablate in the full training scenario using all testing scenes of the DTU dataset [1].

Effect of fine-tuning strategies In this study, we showcase the outcomes of training our occupancy model while fine-tuning different components of GeoNeRF [18]. As illustrated in Table 3, a better occupancy field is learnt when the

Method	Chamfer ↓
w/o tuning encoder	1.29
w tuning encoder	1.15

Table 3: Effect of fine-tuning encoder/feature vol. of GeoNeRF [18]

	w/o $\mathcal{L}_{\nabla o}$	w $\mathcal{L}_{\nabla o}$
NC. ↑	0.64	0.68

Table 4: Normal Consistency

encoder of the pre-trained GeoNeRF [18] model is tuned jointly with the occupancy network. This is because without doing so, the encoder is more suited for learning a density field as is the case for GeoNeRF [18], while we aim to learn an occupancy field.

Effect of Different Loss Components In this study, we showcase the outcomes of various loss components to illustrate their efficacy. As illustrated in the 1st row of Table 5, the novel volumetric rendering weight loss \mathcal{L}_w is vital in learning a sharp and accurate occupancy field. Without it, the chamfer metric degrades, which demonstrates its importance. The 2nd row illustrates the importance of bootstrapping our occupancy with α predictions. Furthermore, we see in Table 6 that the accuracy chamfer metric expectedly suffers more without \mathcal{L}_w as it is directly responsible for sharper surfaces, and less so for the completeness metric.

Method	Chamfer ↓
w/o \mathcal{L}_w	1.22
w/o \mathcal{L}_o	1.17
Full model	1.15

Table 5: Effect of discarding different loss components

	w/o \mathcal{L}_w	w \mathcal{L}_w
Acc. ↓	0.77	0.67
Comp. ↓	1.67	1.64
Overall. ↓	1.22	1.15

Table 6: Ablation of accuracy & completeness

This is also illustrated in Fig. 5, which shows that the loss results in a weight distribution that agrees with one theoretically induced by an occupancy field *i.e.* one that 1-peaks on the surface boundary, as discussed in Unisurf [38]. We also see in the 2nd row that without our bootstrapped occupancy loss \mathcal{L}_o , our surface reconstruction performance is impacted negatively. Finally, we investigate the effect of discarding the normal loss $\mathcal{L}_{\nabla o}$. The normal loss is important for reducing noise in our reconstruction as witnessed in our qualitative results. This smoothing, can however, impact chamfer distance and it can be challenging to strike a perfect balance between these conflicting objectives. To illustrate the importance of normal loss, we provide normal consistency results in Tab. 4 which demonstrates that our normal loss $\mathcal{L}_{\nabla o}$ smoothens the reconstructions and consequently, improves normal consistency between the reconstructed and GT meshes. The whole study demonstrates that all our losses are indispensable to our model’s final performance.

Training time Here, we address the time of training for our model in Table 6. Given a fully-trained GeoNeRF [18], our model is able to finetune its features rapidly. While training SparseNeuS [29] and ReTR [26] fully on a single RTX A6000 GPU takes almost 3 days, we are able to train to adapt features of a pretrained GeoNeRF [18] in approximately 3.5 hrs to learn occupancy fields.

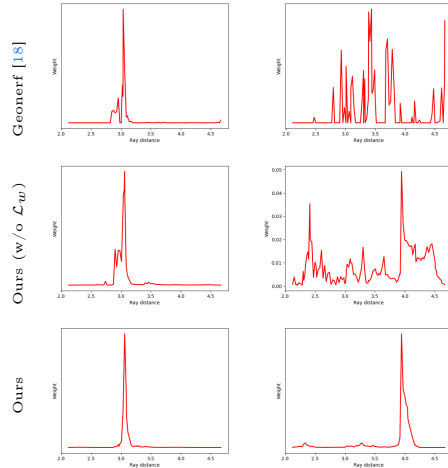


Fig. 5: Comparison between volumetric weight distribution along rays between Geonerf [18] and our method.

5 Limitations

Because of the underlying assumption while learning occupancy in a volumetric framework, our method is specifically designed to represent surfaces that are solid and not transparent. Moreover, the accuracy of the reconstructions is diminished in areas that are infrequently visible in the input images.

6 Conclusions

We introduce a novel approach, GeoTransfer, that rapidly transfers the 3D understanding of generalizable NeRFs to obtain accurate occupancy fields for implicit surface reconstruction. Unlike previous methods that introduced improvements to the encoders of generalizable SDF based reconstruction methods, we discovered that it sufficed to learn a occupancy network in feature space that learns to transform sampling-dependent opacities obtained from the state-of-the-art generalizable NeRF to sampling-independent occupancy fields. This approach also proved to be way faster than training a generalizable sparse 3D surface reconstruction method from scratch. Our method achieves state-of-the-art reconstruction quality for sparse inputs, showcasing its efficacy.

Method	Training duration
SparseNeuS [29]	~3 days
ReTR [26]	~3 days
Ours (re-train)	~3.5 hrs

Fig. 6: Comparison of total training duration.

GeoTransfer: Generalizable Few-Shot Multi-View Reconstruction via Transfer Learning

– Supplementary Material –

Shubhendu Jena, Franck Multon, and Adnane Boukhayma

Inria, Univ. Rennes, CNRS, IRISA, M2S, France

We begin by applying our approach to another baseline model, specifically MVSNeRF [6], and demonstrate significant improvements in its 3D surface reconstruction performance, both quantitatively and qualitatively. Thereafter, we provide additional ablation studies to justify our hyperparameter choices and also show qualitative video comparisons of our reconstruction results to other methods to visually demonstrate the impact of our approach and the corresponding losses. Qualitative comparisons with VolRecon [47] and ReTR [26] for our novel view synthesis results follow, and finally we conclude with some additional experimental details on our evaluation datasets of DTU [1] and Blended-MVS [59].

1 Using MVSNeRF as backbone

To demonstrate the generalizability of our approach, we apply it with MVSNeRF [6] as our backbone and denote the resulting model MVSTransfer. Tuning MVSNeRF [6] to MVSTransfer took us only ~ 3 hrs on a single RTX A6000 GPU. Notice that training MVSNeRF [6] on the same GPU takes at least 3 days. In the following comparative analysis, we are working with the same DTU [1] reconstruction split trained MVSNeRF [6] to ensure fairness. Then, we use the same Sigmoid activated implicit decoder f_o as in our experiments based in GeoNeRF [18], and we leverage the same loss functions to obtain :

$$\mathcal{L} = \sum_{\mathbf{r}} \mathcal{L}_{vol}^o(\mathbf{r}) + \mathcal{L}_{vol}^\sigma(\mathbf{r}) + \lambda \mathcal{L}_{\nabla o}(\mathbf{r}) + \mu \mathcal{L}_w(\mathbf{r}) + \nu \mathcal{L}_o(\mathbf{r}). \quad (1)$$

The losses are weighted exactly as in the original paper and we finetune the decoder for 5400 iterations to obtain a generalizable occupancy network. After running Marching Cubes algorithm [30] with a grid resolution of 400 for extracting each surface mesh by thresholding the occupancy field at 0.5 and computing chamfer metrics with respect to the ground-truth point clouds, we get :

Qualitatively, our approach leads to a globally more accurate learning of the occupancy fields, leading to the extracted mesh to be wrapped more closely around the ground truth meshes (obtained by running screened poisson surface reconstruction on the ground truth point clouds). Some examples are shown in Figure 1 as follows :

As seen above, MVSTransfer (mesh in blue) is closer globally to the ground truth mesh compared to MVSNeRF [6] (mesh in red). This is because of our

Scan	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122
MVSNeRF [6]	2.29	3.76	2.84	1.93	2.79	2.73	1.91	2.51	2.56	2.06	2.01	1.56	1.55	2.24	2.38
MVSTransfer	1.83	3.79	2.5	1.6	2.12	2.51	1.5	1.99	1.99	1.55	1.51	1.53	1.16	1.65	2.0

Table 1: Quantitative comparison on the DTU dataset [1]. The best method is **em-boldened**.

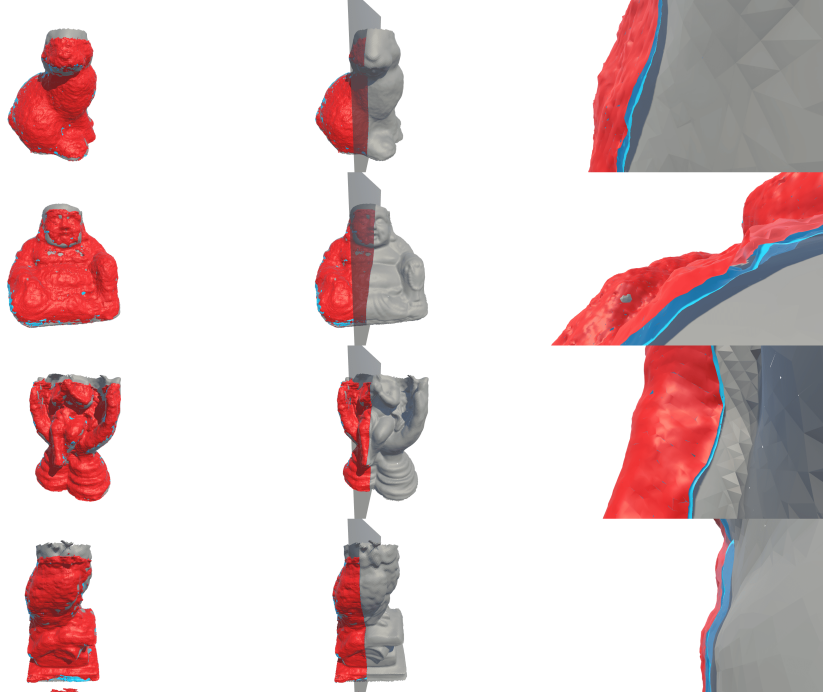


Fig. 1: MVSNeRF [6] (in red) and MVSTransfer (in blue) qualitative evaluation on DTU [1] using 3 source images. **Notice that our meshes (in blue) are closer to the ground truth meshes (in gray) than MVSNeRF [6] (in red)**

losses, particularly the weight rendering loss $\mathcal{L}_w(\mathbf{r})$ which ensures accurate occupancy field estimation in a volumetric rendering framework by reaching a sharp 1-peak at the location of ray-surface intersection.

2 Additional ablation studies

Based on our GeoNeRF experiment [18], We ablate the parameters associated with loss $\mathcal{L}_w(\mathbf{r})$, a_{max} and a_{min} which represent respectively the initial and final widths of the guiding gaussians which 1-peaks at the ray-surface intersection obtained through secant method based root finding. We found the initial value of $a_{max} = 1$ suitable for our training. We ablate the final width a_{min} in Tab.2.

We also ablate the decay parameter β in Tab.3. The latter controls the speed of progression of a in eqn. 14 from a_{max} to a_{min} .

a_{min}	0.002	0.004	0.008
Chamf. ↓	<u>1.16</u>	1.15	1.20

Table 2: Ablation of a_{min}

β	0.0005	0.001	0.002
Chamf. ↓	1.23	1.15	<u>1.17</u>

Table 3: Ablation of β

These studies demonstrate that our chosen a_{min} and β are appropriately chosen to attain accurate generalizable occupancy fields.

3 Additional qualitative comparison on 3D reconstruction

Based on our GeoNeRF experiment [18], we have included some additional video visualizations of our surface reconstructions in the included supplementary material. There are 4 on DTU [1], namely DTU_Scan24.mp4, DTU_Scan55.mp4, DTU_Scan118.mp4 and DTU_Scan122.mp4 and 2 on BlendedMVS [59], namely BMVS_Scan3.mp4, and BMVS_Scan12.mp4.

4 Inference time

We use an occupancy representation as it is akin to the sampling and view dependent NeRF opacity under opaque assumption (*cf.* Unisurf [38]), which facilitates our fast adaptation through transfer learning. Volumetric rendering based inference speed remains similar irrespective of the sdf/occupancy representation. We provide here inference speeds for us and main baselines VolRecon [47] and ReTR [26]. Depth map inference times is about 30s as reported in their respective supplementary sections. We reproduced this on a RTX A6000 and we got 32.4s for us, 31.8s for VolRecon [47] and 37.2s for ReTR [26].

5 Additional qualitative comparison on Novel View Synthesis

In this section, based on our GeoNeRF experiment [18], we present additional qualitative comparisons with VolRecon [47] and ReTR [26] in Figure 2 and 3 to demonstrate the superior performance of our method on the DTU [1] dataset. Our final adapted model preserves the novel-view capabilities of its initial backbone and provides good novel-view extrapolation results compared to the generalizable reconstruction networks (*e.g.* VolRecon [47] and ReTR [26]). We notice that qualitatively, in the sparse 3 input views setting, we are sharper than the competing methods, with lesser artifacts. This demonstrates the robustness of our method on the task of novel-view synthesis, apart from also displaying **state-of-the-art** results on surface reconstruction.

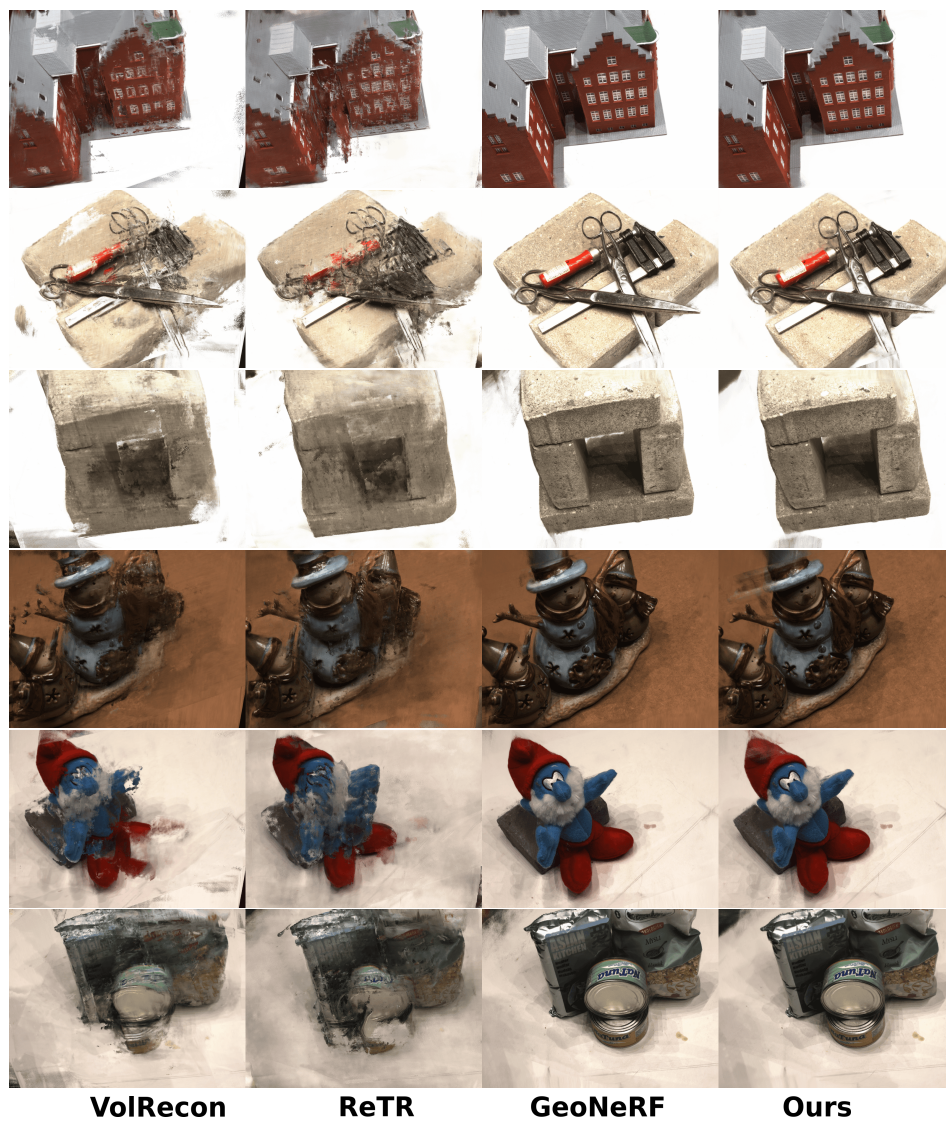


Fig. 2: Novel-View synthesis qualitative evaluation on DTU [1] using 3 source images.



Fig. 3: Novel-View synthesis qualitative evaluation on DTU [1] using 3 source images.

6 Additional Experimental Details

In this work, we evaluated two sets of data: DTU [1] and BlendedMVS [59]. For DTU [1], we follow distinct protocols based on the task’s nature, distinguishing between novel view synthesis and surface reconstruction.

Metrics For the novel view synthesis task involve evaluating PSNR scores, assuming a maximum pixel value of 1 and using the formula $-10 \log_{10}(\text{MSE})$. Additionally, we employ the scikit-image’s API to calculate the Structural Similarity Index (SSIM) score and the pip package lpips, utilizing a learned VGG model for computing the Learned Perceptual Image Patch Similarity (LPIPS) score. In the context of the surface reconstruction task, we gauge Chamfer Distances by comparing predicted meshes with the ground truth point clouds of DTU scans. The evaluation process follows the methodology employed by SparseNeuS, VolRecon, ReTR [26, 29, 47], employing an evaluation script that refines generated meshes using provided object masks. Subsequently, the script evaluates the chamfer distance between sampled points on the generated meshes and the ground truth point cloud, producing distances in both directions before providing an overall average, typically reported in evaluations. Additionally, two sets of 3 different views are used for each scan, and we average the results between the two resulting meshes from each set of images and report it in the comparison as done in previous methods [26, 29, 47].

DTU Dataset The DTU dataset [1] is an extensive multi-view dataset comprising 124 scans featuring various objects. Each scene is composed of 49–64 views with a resolution of 1600×1200 . We adhere to the procedure outlined in [26, 29, 47], training on the same scenes as employed in these methods and then test on the 15 designated test scenes for both the reconstruction and novel view synthesis tasks. The test scan IDs for both novel view synthesis and surface reconstruction are : 24, 37, 40, 55, 63, 192, 65, 69, 83, 97, 105, 106, 110, 114, 118 and 122. For surface reconstruction, for each scan, there are two sets of 3 views with the following IDs used as the input views: set-0: 23, 24 and 33, then set-1: 42, 43 and 44 all scans. We use the training views in half resolution, *i.e.* 800×600 . As for novel-view synthesis, we test for camera views not used during training for a fair evaluation, with the following IDs used as the target and source views for all the scenes - 37 : 39, 36 and 20; 38 : 39, 37 and 40; 39 : 40, 37 and 36.

BlendedMVS Dataset BlendedMVS [59] is a large-scale dataset for generalized multi-view stereo that consists of a variety of 113 scenes including architectures, sculptures and small objects with complex backgrounds. For surface-reconstruction, we use 7 challenging scenes, in accordance with SparseNeuS [29] where each scene has 31–143 images captured at 768×576 . The chosen IDs for the selected scenes are : Scan2 : 67, 29, 59; Scan3 : 1, 0, 2; Scan12 : 2, 8, 50; Scan13: 28, 4, 11; Scan14: 9, 109, 50; Scan22: 4, 3, 5; Scan24: 23, 39, 5. We use the testing views in their original resolution.

Masked evaluation for novel view synthesis In accordance with the findings presented by [36] addressing the bias in background evaluation, we adopt their prescribed approach of masked evaluation for DTU [1]. This involves employing object masks and computing PSNR exclusively within the defined mask. In the case of SSIM and LPIPS, we utilize the masks to superimpose the predicted object-of-interest onto a black background before metric calculations.

References

1. Aanaes, H., Jensen, R.R., Vogiatzis, G., Tola, E., Dahl, A.B.: Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision* **120**(2), 153–168 (2016) [1](#), [2](#), [8](#), [10](#), [11](#), [12](#), [15](#), [16](#), [17](#), [18](#), [19](#), [20](#), [21](#)
2. 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](#)
3. Atzmon, M., Lipman, Y.: Sal: Sign agnostic learning of shapes from raw data. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 2565–2574 (2020) [3](#)
4. Ben-Shabat, Y., Koneputugodage, C.H., Gould, S.: Digs: Divergence guided shape implicit neural representation for unoriented point clouds. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 19323–19332 (2022) [3](#)
5. 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) [4](#)
6. 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) [3](#), [4](#), [8](#), [10](#), [11](#), [15](#), [16](#)
7. 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) [3](#)
8. Darmon, F., Bascle, B., Devaux, J.C., Monasse, P., Aubry, M.: Improving neural implicit surfaces geometry with patch warping. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 6260–6269 (2022) [2](#), [3](#)
9. 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) [3](#)
10. Ding, Y., Yuan, W., Zhu, Q., Zhang, H., Liu, X., Wang, Y., Liu, X.: Transmvsnet: Global context-aware multi-view stereo network with transformers. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 8585–8594 (2022) [1](#)
11. Falcon, W.A.: Pytorch lightning. GitHub **3** (2019) [10](#)
12. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: *CVPR* (2017) [2](#)
13. Gropp, A., Yariv, L., Haim, N., Atzmon, M., Lipman, Y.: Implicit geometric regularization for learning shapes. *arXiv preprint arXiv:2002.10099* (2020) [3](#)
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) [1](#)
15. 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) [3](#)

16. Jena, S., Multon, F., Boukhayma, A.: Neural mesh-based graphics. In: European Conference on Computer Vision. pp. 739–757. Springer (2022) [2](#)
17. Jiang, Y., Ji, D., Han, Z., Zwicker, M.: Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1251–1261 (2020) [3](#)
18. 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) [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [10](#), [11](#), [12](#), [13](#), [14](#), [15](#), [16](#), [17](#)
19. Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: CVPR (2018) [2](#)
20. Kellnhöfer, P., Jebe, L.C., Jones, A., Spicer, R., Pulli, K., Wetzstein, G.: Neural lumigraph rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4287–4297 (2021) [3](#)
21. 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) [2](#)
22. 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. pp. 12912–12921 (2022) [3](#)
23. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014) [10](#)
24. 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) [3](#)
25. 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) [3](#)
26. Liang, Y., He, H., Chen, Y.: Retr: Modeling rendering via transformer for generalizable neural surface reconstruction. *Advances in Neural Information Processing Systems* **36** (2024) [2](#), [4](#), [8](#), [10](#), [11](#), [12](#), [13](#), [14](#), [15](#), [17](#), [20](#)
27. Liu, L., Gu, J., Zaw Lin, K., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. *Advances in Neural Information Processing Systems* **33**, 15651–15663 (2020) [3](#)
28. 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) [3](#), [4](#)
29. Long, X., Lin, C., Wang, P., Komura, T., Wang, W.: Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In: European Conference on Computer Vision. pp. 210–227. Springer (2022) [2](#), [4](#), [8](#), [10](#), [11](#), [13](#), [14](#), [20](#)
30. Lorensen, W.E., Cline, H.E.: Marching cubes: A high resolution 3d surface construction algorithm. In: *Seminal graphics: pioneering efforts that shaped the field*, pp. 347–353 (1998) [4](#), [10](#), [15](#)
31. Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* (2016) [10](#)
32. Max, N.: Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics* **1**(2), 99–108 (1995) [5](#)
33. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4460–4470 (2019) [3](#)

34. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* **38**(4), 1–14 (2019) [3](#), [8](#)
35. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* **65**(1), 99–106 (2021) [1](#), [3](#), [4](#), [6](#)
36. 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) [2](#), [3](#), [21](#)
37. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3504–3515 (2020) [2](#), [3](#), [7](#)
38. 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) [2](#), [3](#), [4](#), [7](#), [8](#), [10](#), [13](#), [17](#)
39. Ouasfi, A., Boukhayma, A.: Few ‘zero level set’-shot learning of shape signed distance functions in feature space. In: *European Conference on Computer Vision*. pp. 561–578. Springer (2022) [4](#)
40. 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> [3](#)
41. Ouasfi, A., Boukhayma, A.: Mixing-denoising generalizable occupancy networks. In: *2024 International Conference on 3D Vision (3DV)*. pp. 1103–1114. IEEE (2024) [4](#)
42. Ouasfi, A., Boukhayma, A.: Robustifying generalizable implicit shape networks with a tunable non-parametric model. *Advances in Neural Information Processing Systems* **36** (2024) [4](#)
43. Ouasfi, A., Boukhayma, A.: Unsupervised occupancy learning from sparse point cloud. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 21729–21739 (June 2024) [3](#)
44. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 165–174 (2019) [2](#), [3](#)
45. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019) [7](#), [10](#)
46. Peng, S., Niemeyer, M., Mescheder, L., Pollefeys, M., Geiger, A.: Convolutional occupancy networks. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III* 16. pp. 523–540. Springer (2020) [4](#)
47. Ren, Y., Zhang, T., Pollefeys, M., Süsstrunk, S., Wang, F.: Volrecon: Volume rendering of signed ray distance functions for generalizable multi-view reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 16685–16695 (2023) [2](#), [4](#), [8](#), [10](#), [11](#), [12](#), [15](#), [17](#), [20](#)

48. 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) [8](#), [10](#), [11](#)
49. 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) [3](#)
50. Wang, F., Galliani, S., Vogel, C., Speciale, P., Pollefeys, M.: Patchmatchnet: Learned multi-view patchmatch stereo. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 14194–14203 (2021) [1](#)
51. 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) [2](#)
52. 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) [2](#), [3](#), [7](#), [8](#), [10](#)
53. 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) [3](#), [4](#), [8](#), [10](#)
54. 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) [3](#)
55. Wynn, J., Turmukhambetov, D.: Diffuserf: 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) [3](#)
56. 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) [2](#)
57. Yang, J., Mao, W., Alvarez, J.M., Liu, M.: Cost volume pyramid based depth inference for multi-view stereo. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4877–4886 (2020) [1](#)
58. Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: Proceedings of the European conference on computer vision (ECCV). pp. 767–783 (2018) [1](#), [8](#), [10](#), [11](#)
59. Yao, Y., Luo, Z., Li, S., Zhang, J., Ren, Y., Zhou, L., Fang, T., Quan, L.: Blended-mvs: 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) [1](#), [8](#), [10](#), [11](#), [15](#), [17](#), [20](#)
60. Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. Advances in Neural Information Processing Systems **34**, 4805–4815 (2021) [2](#), [3](#), [8](#), [10](#)
61. 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) [3](#), [8](#), [10](#)
62. Younes, M., Ouasfi, A., Boukhayma, A.: Sparsecraft: Few-shot neural reconstruction through stereopsis guided geometric linearization. In: ECCV (2024) [3](#)
63. 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) [3](#), [4](#), [8](#), [10](#)

- 64. Yu, Z., Peng, S., Niemeyer, M., Sattler, T., Geiger, A.: Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems* **35**, 25018–25032 (2022) [3](#)