

MegaSurf: Scalable Large Scene Neural Surface Reconstruction

Yusen Wang*
wangyusen@whu.edu.cn
Wuhan University
School of Computer Science
Wuhan, Hubei, China

Wenxiao Zhang
wenxxiao.zhang@gmail.com
University of Science and Technology of China
School of Information Science and Technology
Hefei, Anhui, China

Kaixuan Zhou*
zhoukaixuan2@huawei.com
Huawei Technologies
Riemann Lab
Wuhan, Hubei, China

Chunxia Xiao[†]
cxxiao@whu.edu.cn
Wuhan University
School of Computer Science
Wuhan, Hubei, China

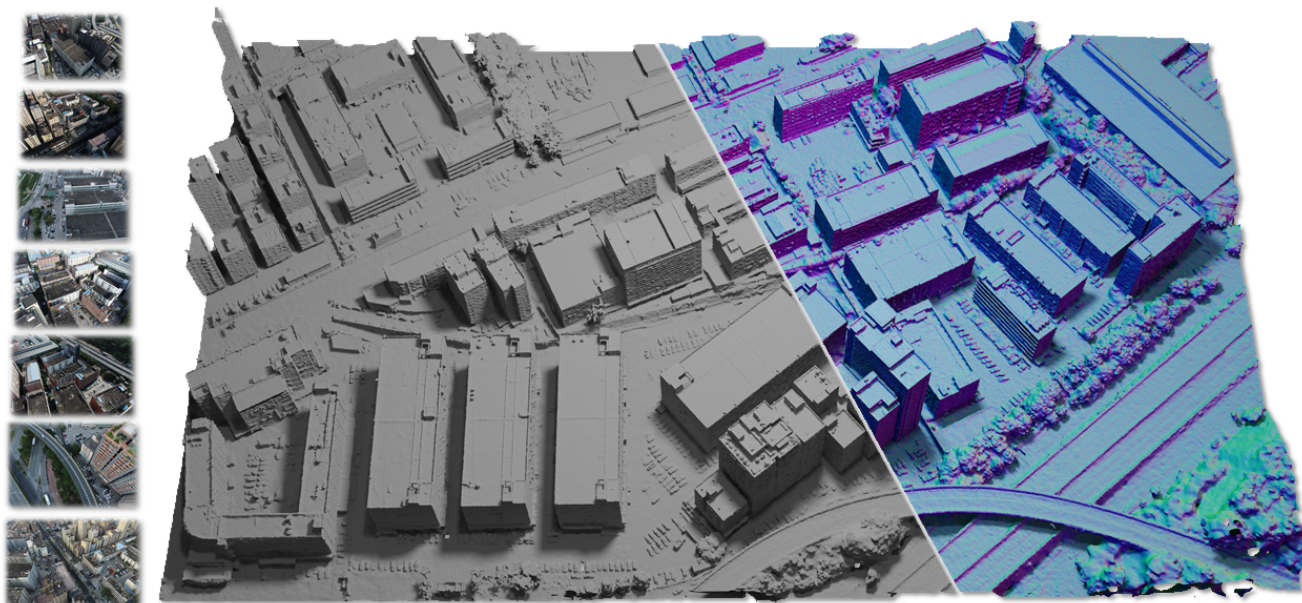


Figure 1: MegaSurf is designed to reconstruct large-scale scenes from extensive images captured by drones. It has both the robustness of the stereo matching and the high-fidelity details of the rendering-based reconstruction methods.

Abstract

We present MegaSurf, a Neural Surface Reconstruction (NSR) framework to reconstruct 3D models of large scenes from aerial images. Many methods utilize geometry cues to overcome the shape-radiance ambiguity, which would produce large geometric errors. In addition, directly using inevitable imprecise geometric cues would lead to degradation in the reconstruction results, especially on large-scale scenes. To address this phenomenon, we propose a Learnable Geometric Guider (LG Guider) to learn a sampling field from reliable geometric cues. The LG Guider decides which position should fit the input radiance and can be continuously refined by rendering loss. Our MegaSurf uses a Divide-and-Conquer training strategy to address the synchronization issue between the Guider and the lagging NSR’s radiance field. This strategy enables the Guider to

transmit the information it carried to the radiance field without being disrupted by the gradients back-propagated from the lagging rendering loss at the early stage of training. Furthermore, we propose a Fast PatchMatch MVS module to derive the geometric cues in the planer regions that help overcome ambiguity. Experiments on several aerial datasets show that MegaSurf can overcome ambiguity while preserving high-fidelity details. Compared to SOTA methods, MegaSurf achieves superior reconstruction accuracy of large scenes and boosts the acquisition of geometric cues more than four times.

CCS Concepts

• **Computing methodologies** → **Reconstruction.**

Keywords

Neural Surface Reconstruction, Large Scale Scenes, Multiview Reconstruction

*Both authors contributed equally to this research.

[†]Corresponding author

1 Introduction

Recently, Neural Surface Reconstruction (NSR), derived from neural radiance field (NeRF) [21, 30, 44, 50], not only excels in high-fidelity novel view synthesis, but also enables accurate geometric acquisition. Accurate 3D surface model is an essential and basic element to create immersive experiences in the game engines and VR experiences. Although NSR has achieved good results in small-scale scenes, there is limited research on its effectiveness in large-scale scenes. Besides, there are many works [11, 28, 29, 36, 39, 41] on the novel view synthesis of large scenes, but little research on the 3D reconstruction directly using images without the aid of LiDAR [6, 23]. Utilizing drones for image acquisition and employing NSR technology can efficiently digitize and vividly recreate cities and historical sites for preservation, while also supporting wide-spreading of AR/VR applications.

However, NSR often encounters geometric errors due to shape-radiance ambiguity, as NSR uses rendering loss to optimize geometries for SDF network [3] implicitly. This problem becomes worse when NSR needs to render aerial images captured for large and complex scenes. Therefore, existing works [3, 34, 49] introduce geometric cues from multi-view stereo into NSR, imposing additional geometric constraints on the rendering, thereby improving the accuracy of the NSR methods.

Methods like MonoSDF [34, 49] add a geometric loss by rendering depth from network to compare with the geometric cues provided by depth estimation. Some other methods like GeoNeuS [3, 31] use a multiview photometric consistency loss derived from the implicit surface as the geometric loss without explicitly deriving geometries from MVS. However, geometric errors and noise persist in the radiance field due to strong constraints caused by geometry loss. This prevents the implicit surface from fitting the real geometry accurately and leads to the degradation of details. NerfingMVS [37] employs confidence of geometric cues to define the sampling range around the prior to deal with noisy geometric cues. However, the confidence of geometric cues is difficult to assess, and manually set the threshold to the sampling ranges is too rigid to be applied to different and variable datasets.

We propose a Learnable Geometric Guider (LG Guider) which firstly distill the geometric cues to the sampling network to avoid sampling on the ambiguous regions, and also can be continuously refined by rendering to overcome the missing details due to noises of geometric cues. If the LG Guider is used to guide an unoptimized radiance field directly, the learned geometric information carried by the Guider will be damaged and causing ambiguity again (Fig. 3). Therefore, we propose the Divide-and-Conquer training strategy as shown in Fig. 2. Firstly, we train the LG Guider with geometry net with geometry cues, to distill the prior geometric cues to prior knowledge of sampling and SDF field. Then we freeze the LG Guider, and train render net. The purpose is to use the distilled sampling retrain the radiance field falling into ambiguous regions. Finally we train the full network, to refine the sampling and geometry by rendering loss for recovering geometric details from noisy geometric cues. To be noticed, geometry cues is only introduced in the first stage to avoid its continuous noise effects to the final results.

In addition, we find that ambiguities in NSR often occur in the large planar geometries in the region of low texture and shadows, while the complex geometries often easier to be reconstructed due to their rich and distinct color. We propose a fast PatchMatch MVS module to efficiently reconstruct the large planar geometries. A novel local propagation strategy is designed which progressively propagate geometries with similar plane with the SFM points, with only one step of PatchMatch operation performed per pixel to speed up the MVS process.

In summary, our main contributions are the following:

- We introduce a Learnable Geometric Guider to distill the geometric cues to overcome shape-radiance ambiguities and can be continuously refined by rendering to recover details from geometric noises.
- We propose a Divide-and-Conquer training strategy to improve the guidance of learning of shape and radiance field using the Learnable Geometric Guider.
- We present a fast MVS module to efficiently obtain high confidence planar geometric priors over 4× improvement in speed where large shape-radiance ambiguities often occur.
- On the several aerial photography datasets, our algorithm achieved the best results of quantitative and qualitative results. To our knowledge, we are the first to extend accurate NSR to large scale aerial scene.

2 Related Work

Multiview stereo matching. Multiview Stereo (MVS) [26] aims to recover 3D geometric model of the real scene from input images. The key idea of image based multiview reconstruction is photo-consistency matching [4, 15, 25, 43]. However, the performance of local photo-consistency matching is easily reduced in regions with low textures, shadows, and non-Lambertian materials. Therefore, several global matching aggregation methods are applied to improve the quality, including semi-global optimization [7], PatchMatch [25], and 3D convolution regularization [45]. Even though the learning-based MVS methods [5, 10, 14, 17, 20, 40, 45, 52] show their advantages of reconstruction in difficult regions, their application on large-scale aerial datasets is limited due to the lack of various 3D training datasets, which are often expensive to acquire. Patchmatch-based MVS methods [25, 42], with their efficient parallelization structure and robust performance, are more suitable and already widely applied for large-scale scene reconstruction. However, common PatchMatch MVS requires performing PatchMatch operations several times through all pixels globally from random initialization. These intensive computation especially on large scale datasets introduce an unnegligible overhead when using them as geometric cues for NSR.

Neural surface reconstruction. Recently, rendering-based neural surface reconstruction methods [12, 19, 27, 32, 46, 47] have become a promising way to promote the development of 3D reconstruction due to their high-quality reconstruction results, especially for fine structures [13, 35] and training speed [24, 33, 38]. The multi-resolution hash encoding [22] provides a compact high-resolution feature representation which shows its potential for high-fidelity reconstruction for large scenes. Li et al [13] introduce a progressive

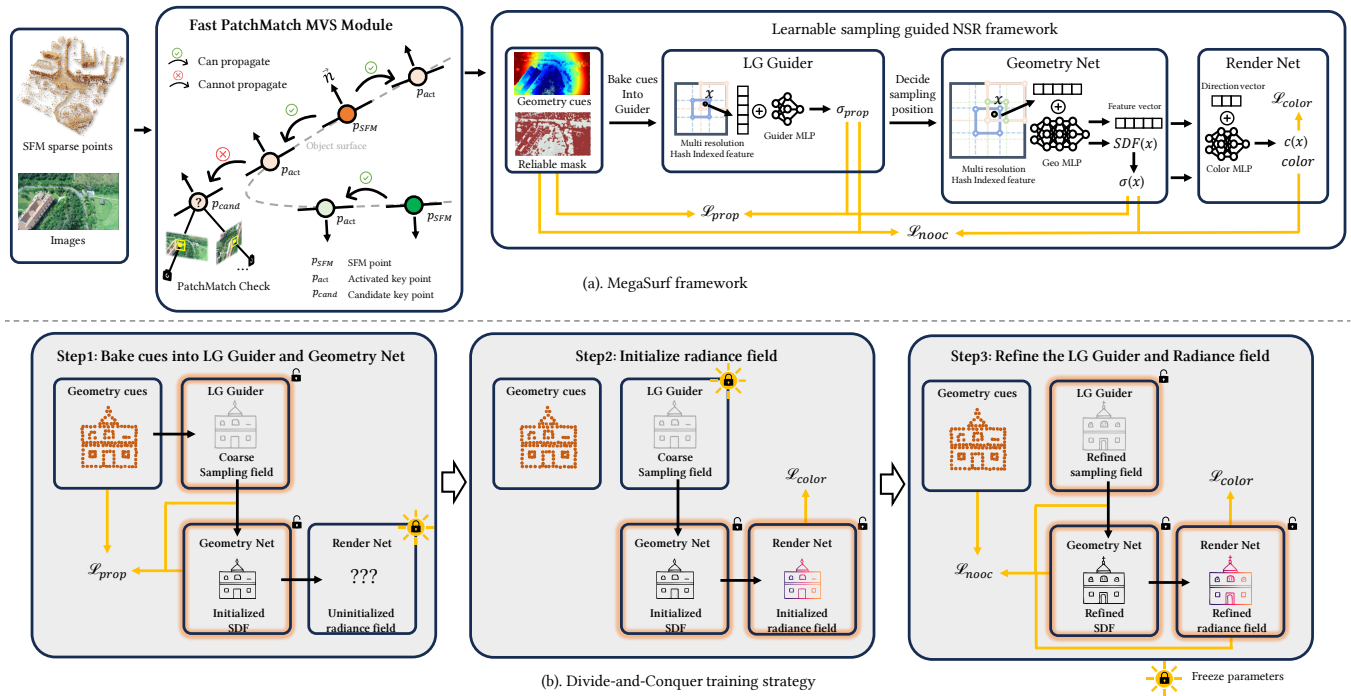


Figure 2: Method overview. (a) We propose a Fast PatchMatch MVS Module (Section. 3.4) to rapidly propagate SFM points to obtain high-confidence geometry cues. Then, we use these cues to train our Learnable Geometric Guider (LG Guider, Section 3.2). The LG Guider which position should be used to fit the input radiance and can be continuously refined by rendering loss. To address the synchronization issue between the Guider and the lagging radiance field, we propose a three steps Divide-and-Conquer training strategy (Section. 3.3). This strategy enables the Guider to efficiently guide the radiance field training while preserving its learned geometric information from being disrupted by the rendering. (b) The detail of our Divide-and-Conquer training strategy. We bake the geometry cues into LG Guider and Geometry net in Step 1, then freeze the LG Guider parameters and initialize the whole radiance field in Step 2. In Step 3, we use rendering loss to refine the radiance field and our LG Guider and propose L_{prop} to preserve the geometry information the Guider carries from being impaired.

training strategy on the multi-resolution hash encoding representation, and a numerical calculation of normals, firstly extending the high reconstruction accuracy to large outdoor scenes. However, on large-scale aerial scenes, large geometrical errors often occur due to the more server shape-radiance ambiguity in the complicate large scenes as show by Neuralangelo [13] results in Figure 6.

Neural surface reconstruction with geometry cues. Many works incorporate geometry cues into NSR reconstruction to address the ambiguity problem. Several of these utilize the geometry cues as a geometry loss to ensure that the geometry reconstructed from NSR are consistent with the geometric cues. However, noisy geometry cues persistently contribute to the loss, resulting in over-smooth effects on the detailed structures. To avoid the intensive computation of global optimization [7, 25] of MVS, some other works directly use the photo consistency measurement, Normalized Cross-Correlation (NCC) as geometry cues. However, NCC, a highly localized geometric measurement, often fails to give reliable geometry in the ambiguous areas, resulting in a worse reconstruction in the large scenes(see Supplementary Materials for details).

Another approach retrains the sampling points around the geometric cues to deal with noisy geometric cues. Wei et al employ the confidence of geometric cues to define the sampling range around the prior to deal with noises. However, the confidence of geometric cues is difficult to assess and manually set the threshold to the sampling ranges cannot be applied to different and variable datasets.

3 Method

As shown in Figure 2, MegaSurf proposes a Fast PatchMatch MVS Module to efficiently obtain the geometry cues (Section 3.4). Then, we propose a Learnable Geometric Guider (Section 3.2) to learn these reliable geometry cues. Next, MegaSurf employs a Divide-and-Conquer training strategy (Section 3.3) to train the radiance field.

3.1 Preliminary

Neural radiance field. NeRF [21] represents a complex 3D scene as a learned function that maps each 3D point and corresponding ray direction to a color and density. It integrates the color of sampled

points along the ray to render each pixel:

$$C(r) = \sum_i \omega_i \mathbf{c}_i, \quad \omega_i = T_i \alpha_i, \quad \alpha_i = 1 - \exp(-\sigma_i \delta_i), \quad (1)$$

where α_i is the opacity of the i th sample point along the ray, σ_i is the corresponding density, which is also the learned function’s output. $\delta_i = t_i - t_{i-1}$, is the distance between two sample points, t is the distance to the ray center. $T_i = \prod_{j=1}^{i-1} (1 - \sigma_j)$ is the accumulated transmittance. As the geometry of NeRF is represented by density, extracting surfaces from densities often leads to noisy results.

Neural surface reconstruction. Most rendering based NSR methods take NeRF as the backbone and use signed distance function (SDF) as the geometric representation instead of density in NeRFs. The surface can be represented by the zero-level set of the SDF, $S = \{\mathbf{x} : f(\mathbf{x}) = 0\}$, where x is a 3D position. To use volume rendering, VolSDF [46] defines the volume density function τ to map the signed distance $f(x)$ to volume density σ :

$$\tau(\mathbf{x}) = \beta^{-1} \Psi_\beta(f(\mathbf{x})), \quad (2)$$

where $\beta > 0$ is a scheduling parameters and approaches 0 during optimization, $\tau(\mathbf{x})$ is the cumulative distribution function (CDF) of the zero-mean Laplace distribution with scale β . Manually controlling the β allows different reconstructed cases to have the same β , so that the surface details of different cases are consistent.

Neuralangelo. Recently, multi-resolution hash encoding proposed by Muller et al. [22] is a compact feature representation that can represent large-scale scenes in unprecedented detail. Neuralangelo [13] designs a coarse-to-fine optimization scheme to reconstruct the surfaces with progressive levels of detail:

$$\gamma_l = [F_0, F_1, \dots, F_{start+l}], \quad l_{start} < l < l_{max}, \quad (3)$$

where γ represents the features from hash grids, F is the features of each level of hash grid, and the coarse to fine resolution spans from level l_{start} to level l_{max} . Another important contribution is the design of a numerical gradient computation to distribute the back-propagation updates to wider neighboring hash grids to improve the smoothness of surface reconstruction:

$$\nabla_x f(x) = \frac{f(\gamma(x + \epsilon_x)) - f(\gamma(x - \epsilon_x))}{2\epsilon}, \quad (4)$$

where ϵ is the step size away from x for sampling points to calculate gradient numerically.

However, when applying it to large-scale aerial datasets, severe shape radiance often happens in the areas of heavy shadows, low textures, and illumination variations.

3.2 Learnable sampling guided NSR

Our learnable geometric guider borrows the sampling proposal network to distill the geometric clues to restrain samplings in the ambiguity areas. The proposal net adopts a two-level, (coarse level: $Prop_0$ and fine level: $Prop_1$), coarse to fine hierarchical sampling procedure [21]. Each level consists of a small multi-resolution hash grid and a tiny MLP to learn the importance of sampling to propose informative samples to subsequent geometry and render net to learn the SDF and radiance field as shown in Figure 2. A naive solution is using geometric cues to train the LG Guider while simultaneously training the whole network. This approach is similar to [37] which constrains the sampling during training, but our

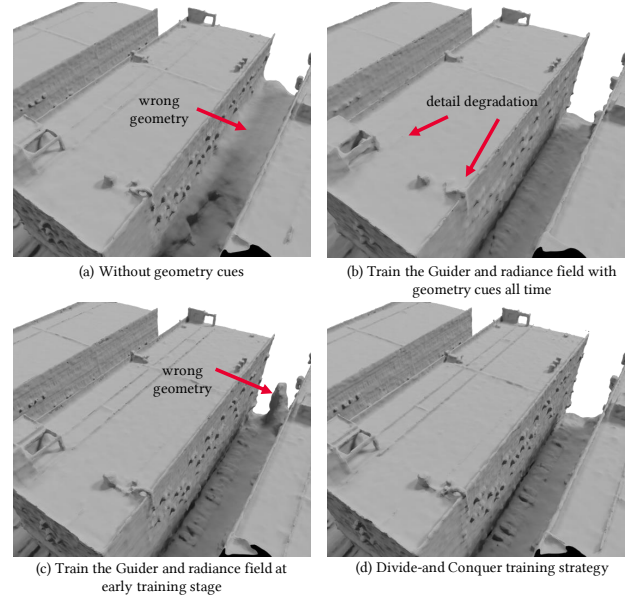


Figure 3: The illustration of the impact of different training strategies using geometry cues. (a) Training without geometry cues. (b) Train the LG Guider and radiance field simultaneously using geometry cues all the time. (c) Only use geometry cues to train the LG Guider and the radiance field at early training stage. (d) Our Divide-and-Conquer training strategy showing in Figure 2. For ablation studies on all settings, please refer to the supplementary material.

sampler is learnable. As the geometric noises affect the training process all the time, this solution would also inevitably introduce noisy geometry cues to the final reconstruction, leading to missing details as shown in Figure 3b. Another approach is to bake the geometric cues by training LG Guider and the whole network at the beginning of the training and then letting the network refine the noisy geometric guider without using geometric cues. However, when introducing the render net into the baking step, the ambiguity problem of rendering will directly affect the baking results, leading to some remaining geometric errors in the final reconstruction, as shown in Figure 3c.

3.3 Divide-and-Conquer training strategy

We propose a Divide-and-Conquer training strategy to distill the geometric information into the LG Guider and enable it to specify specific positions in the radiance field for optimization. The Divide-and-Conquer training strategy consists of three steps: 1) Baking geometry cues into LG Guider, 2) Initializing the radiance field, 3) Refining training.

Step1: Bake cues into LG Guider. In baking stage, we train the LG Guider with geometry net with geometry cues, to distill the prior geometric cues to prior knowledge of sampling and SDF field while leaving render net untrained as shown in Figure 2 bottom left. The reason for not training the render net at this stage is that when rendering loss is introduced, the introduced ambiguity problem can

affect the distillation of prior geometric knowledge to both the LG Guider and the geometry net. This reduces the effectiveness of LG Guider to resolving ambiguities as shown in bottom left in Figure 3

To be specific, we maximize the sampling weights ω^h given by coarse level $Prop_0$ and fine level $Prop_1$ and the radiance field weight ω^{geo} given by geometry net within the range $[t_{prior} - \epsilon, t_{prior} + \epsilon]$ around the geometry cues t_{prior} :

$$L_{prop} = 1 - \sum_{i \in \Lambda} (\omega_i^h + \omega_i^{geo}), \quad (5)$$

$$\Lambda : \{i : t_{prior} - \epsilon < t_i < t_{prior} + \epsilon\}, h \in Prop_0, Prop_1,$$

where t_{prior} is the distance between the camera center to the 3D point corresponding to the depth cue. The computation of ω_i^h follows the Eqn. 1, we replace the geometry net output σ with the LG Guider output σ_{prop} .

We further add a Curvature loss to improve the smoothness of sampling field and the geometry to address the noise and incompleteness of the geometric cues:

$$\mathcal{L}_{curv} = \frac{1}{N} \sum_{i=1}^N |\nabla^2 f(x_i)|, \quad (6)$$

The overall loss of step 1 is:

$$L_{step1} = L_{prop} + L_{curv}. \quad (7)$$

In this way, step 1 completes the training of the LG Guider and also finishes the initialization of the geometry net. Next, we need to initialize the entire radiance field by adding color information to the geometry represented by the geometry net.

Step2: Initialize radiance field. Since our geometry net and render net share a multi-resolution hash grid to provide hash indexed feature vectors, modifications to the render net will impact the geometry net. However, the structure of the LG Guider is independent of the geometry net and render net. To prevent the learned information of the LG Guider from being compromised by the uninitialized render net, we freeze the parameters of the LG Guider.

Since the parameters of the LG Guider are fixed, the sampling positions outputted by the LG Guider are also fixed. The radiance field will only prioritize using these positions to fit the input colors, which will help overcome shape-radiance ambiguity.

The step 2 training loss is defined as:

$$L_{step2} = L_{color} + L_{curv} + L_{eikonal}. \quad (8)$$

We take rendering loss L_{color} as primary loss, and take Curvature loss L_{curv} and Eikonal loss $L_{eikonal}$ as regularization terms.

Step3: Refine the LG Guider and radiance field. In this step, we unfreeze all parameters for training, aiming to use rendering loss to refine the LG Guider which is affected by prior noisy geometry. During this process, we further employ the prior geometry cues to avoid the rendering to step back into ambiguity regions, hence we introduce Non occupancy loss L_{nocc} :

$$L_{nocc} = \left\| \sum_{i \in \Gamma} \omega_i c_i \right\|_1, \quad (9)$$

$$\Gamma : \{i : t_i < t_{prior} - \epsilon\},$$

where ω and c is given by geometry net and render net. L_{nocc} is used to ensure that no new surfaces appear between the camera

center and the surfaces corresponding to reliable cues. Thus that the accumulated color should be nearly black color and L_{nocc} should be close to 0. As the LG Guider decides the radiance field’s sampling positions, the rendering loss can be backpropagated to LG Guider, which makes the sampling more precise. We add a Non occupancy loss into the loss of step 2:

$$L_{step3} = L_{color} + L_{nocc} + L_{curv} + L_{eikonal}. \quad (10)$$

After training is completed, we utilize Marching Cube [18] to extract the zero level set from the signed distance function (SDF) represented by the geometry net as the final reconstructed mesh.

3.4 Fast PatchMatch MVS Module

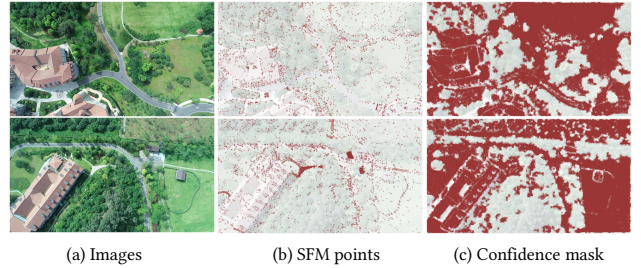


Figure 4: The illustration of the high-confidence region acquired by our Fast PatchMatch MVS module. (a) The input images. (b) Sparse SFM points. (c) The high-confidence position which we used as the geometric prior during our NSR training.

Preliminary of heavy PatchMatch MVS module. Commonly used PatchMatch MVS module starts from randomly initializing geometry on each pixel, and every pixel uses PatchMatch optimization to select its best geometric candidate with the smallest photo-consistency loss E_{NCC} from all the candidates propagated from its neighboring pixels [25, 42]. In a nutshell, PatchMatch operation is to choose the best geometric candidate propagated from neighborhoods for each pixel. Every pixel will continuously update its geometry through PatchMatch until it receives its accurate geometry. Due to the random initialization, pixels often require several (4 times in [42]) global PatchMatch optimizations to get the accurate candidate to converge, which are the major computation cost contribute to MVS.

Fast local propagation from SFM points. Instead, we start from high confident SFM points in each image as activate key points p_{act} to propagate the information to surrounding neighbors. We randomly select eight neighboring pixels for each p_{act} within a 11×11 pixel area as candidate key points p_{cand} . Next, we perform PatchMatch operation on the p_{cand} . The p_{cand} become a new p_{act} when they satisfy that the distance of the p_{cand} to the corresponding p_{act} is less than the given reconstruction accuracy.

In this way, if the p_{cand} is in the similar plane with the p_{act} , it immediately receive the its accurate candidate geometry from p_{act} which will be most likely selected from one-step PatchMatch operation with a minimal photo-consistency loss comparing to other neighboring geometric candidates.

When the activated key point is determined, we design a skip propagation strategy to further propagation by generating a neighbor mask from the activated key as shown in Figure 5. No new key point would be sampled within this mask. This is to mitigate the incorrect propagation to the outside of the plane across the boundary. When no p_{act} exists, we perform the PatchMatch operation for all pixels that are not sampled.

Our propagation strategy ensures every pixel to perform PatchMatch operation once to speed up MVS to more than 4 times. As shown in Figure 4, Our reconstruction aims to reconstruction high confident large planar geometries from SFM points where large shape-radiance ambiguities more likely occur, and leaves non-planar geometries, such as trees and fine details, where NSR methods can reconstruct better.

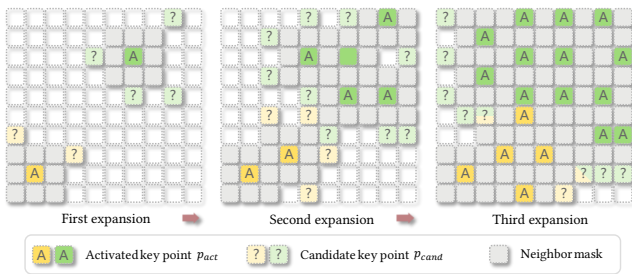


Figure 5: The propagation strategy of our Fast PatchMatch MVS module. The high-confidence geometric information is progressively propagated to its surrounding area.

4 Experiments

4.1 Experimental Setup

Baselines. Our experiments are conducted on UrbanScene3D [16], Mill19 [29] and Songshanhu which is collected by our drone. Their areas are between $60000m^2$ ($300m \times 200m$) and $150000m^2$ ($300m \times 500m$). We divided the whole scenes into several blocks and each block covers a $150m \times 150m$ ground region. We compare MegaSurf with ACMH [42], a traditional reconstruction method, and two NSR methods: Bakedangelo [48] and Monoangelo. Bakedangelo combines BakedSDF [47] with Neuralangelo [13] settings and has a better background modeling, which is more efficient than Neuralangelo. We migrate the key ideas of MonoSDF [49] to Bakedangelo which called Monoangelo, as the results obtained by MonoSDF are generally oversmooth.

We train MegaSurf for 200k iterations per block (step1: 10k, step2 10k, step3: 180k). The memory consumption is about 22G. The efficiency is basically the same as Bakedangelo [48]. The weights of Curvature loss, Eikonal loss, and L_{nooc} are all $1e-3$; the others are all 1. After NSR training, we extract the mesh from the SDF by Marching Cube [18]. We compared the reconstruction results of SciArt and Polytech with the LiDAR ground truth following the official evaluation protocol.

4.2 Comparisons

We developed our Fast PatchMatch MVS module on ACMH software [42], which claims the better quality, and three time speed than another popular open source software, COLMAP [25]. We project the high-confidence geometric cues obtained by our Fast PatchMatch MVS module to the 3D space to form a point cloud and compare it with ACMH.

Table 1: Quantitative results of generating the priors of our Fast PatchMatch MVS module vs ACMH.

Method	$Acc_{50} / Comp_{50} / Overall_{50} \downarrow$	$Acc_{95} / Comp_{95} / Overall_{95} \downarrow$	PM Time \downarrow
SciArt			
ACMH	0.1566 / 0.1432 / 0.1499	0.2035 / 0.3663 / 0.2849	7274s
Ours	0.1629 / 0.1320 / 0.1475	0.2010 / 0.3832 / 0.2921	1357s
Polytech			
ACMH	0.1021 / 0.1043 / 0.1032	0.1701 / 0.2300 / 0.2000	7641s
Ours	0.1227 / 0.1218 / 0.1222	0.1937 / 0.2704 / 0.2320	1460s

Table 2: Quantitative evaluation of reconstruction with existing methods on the UrbanScene3D dataset. MegaSurf achieves the best surface reconstruction performance.

Method	CD \downarrow	$Acc_{95} \downarrow$	$Comp_{95} \downarrow$	$Overall_{95} \downarrow$
SciArt				
ACMH	1.1675	0.2958	0.5136	0.4047
Bakedangelo	1.3938	0.3319	0.5813	0.4566
Monoangelo	1.4142	0.3778	0.6152	0.4965
Ours	<u>1.0574</u>	0.2990	<u>0.4138</u>	<u>0.3564</u>
Polytech				
ACMH	0.6913	<u>0.1588</u>	0.2499	0.2044
Bakedangelo	1.1029	0.2989	0.3969	0.3479
Monoangelo	0.7414	0.1810	0.2472	0.2141
Ours	<u>0.6593</u>	0.1763	<u>0.2086</u>	<u>0.1925</u>

We report quality evaluation results for the top 50% accuracy and 95% accuracy points to reduce the influence of noise. Table 1 shows that the reconstruction accuracy of our module is comparable to ACMH, and the speed of the Patchmatch stage is 4 times faster. This matches the configuration of ACMH, which applies four times PatchMatch global sweeps on each pixel. Furthermore, ACMH requires a depth fusion step to filter noisy geometries for the final geometric cues. This step is extremely slow when a large number of images are applied due to their naive implementation, which is not counted in our table. Note that we do not need this fusion step and can also get comparable geometries with reliable masks.

We provide qualitative and quantitative comparisons to evaluate the performance of our method. Fig 6 and Table 2 shows the results respectively. We achieved the best results in terms of the Chamfer Distance (CD) and Overall score (Overall).

Bakedangelo can generate realistic details but suffers inherent shape-radiance ambiguity due to the lack of geometric constraints, often leading to incorrect geometry. Traditional methods such as ACMH are stable in large scene reconstruction. However, due to the large amount of noise in point clouds, the triangulation may incorrectly connect the points and cause over-smoothing. Monoangelo

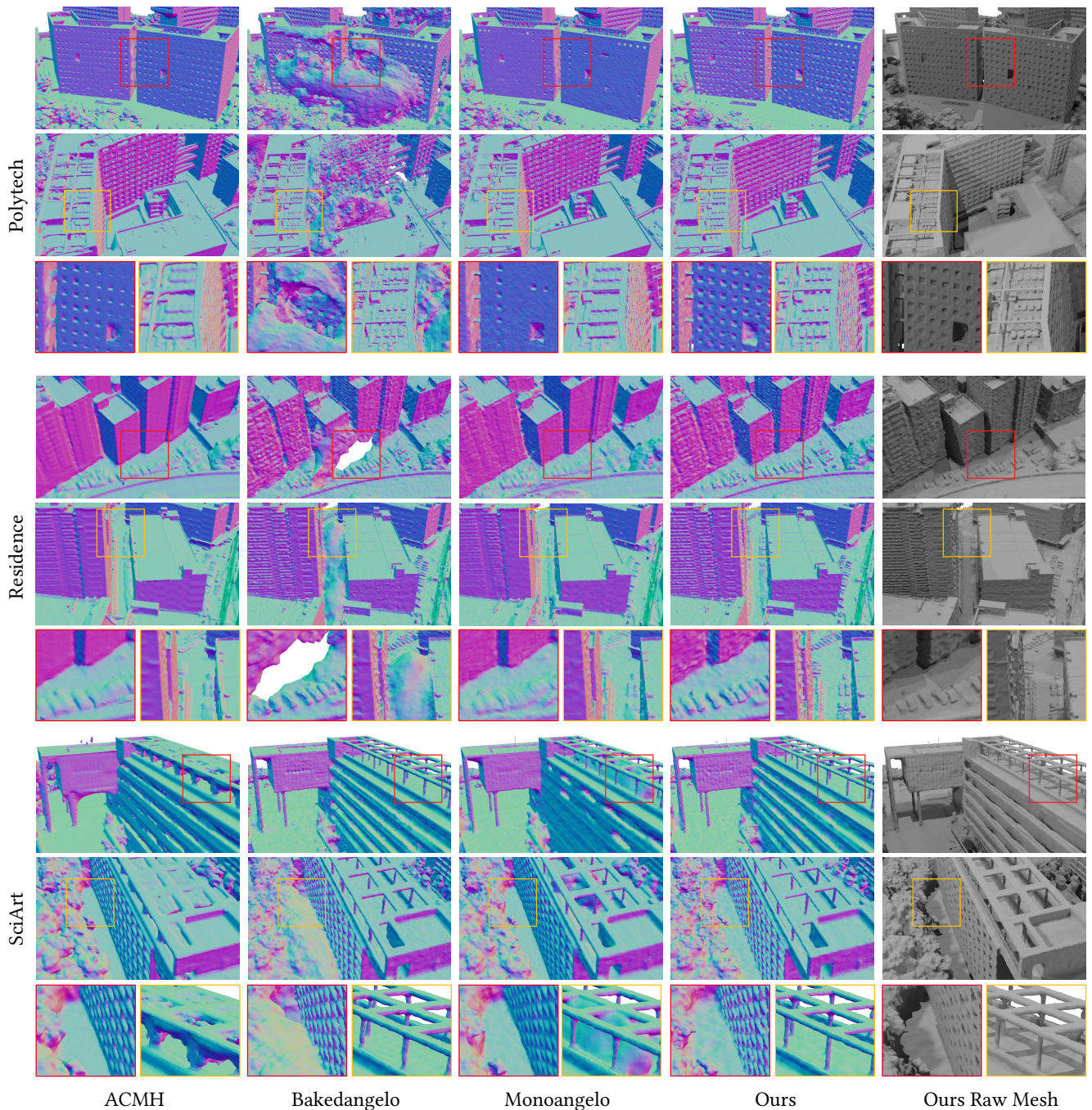


Figure 6: Qualitative results on the UrbanScene3D dataset. MegaSurf both have the robustness to the severe shape-radiance ambiguity and preserve high-fidelity details. The first four columns show the Normal of the corresponding mesh. (Polytech: 12 blocks, Residence: 16 blocks, SciArt: 6 blocks)

takes depth priors as a regular term to guide the NSR optimization. The depth provided by MVS can help Monoangelo overcome shape-radiance ambiguity, but the noise in priors makes it difficult to reconstruct the fine geometric details. Our MegaSurf utilizes

the LG Guider to learn accurate geometric knowledge and prioritize fitting input colors to guide certain regions of the radiance field, thereby overcoming the shape-radiance ambiguity present in Bakedangelo. Additionally, since MegaSurf does not directly employ inaccurate geometry loss and the LG Guider can continuously

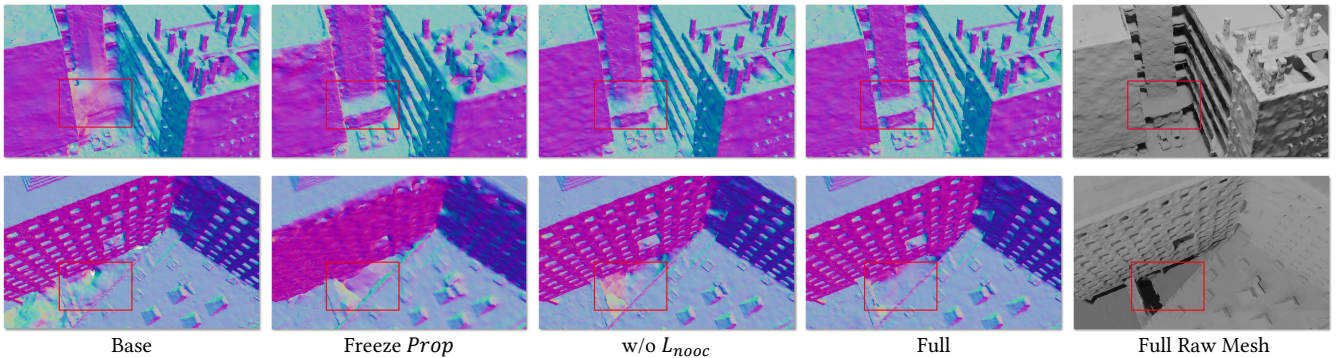


Figure 7: Visualization results of the ablation study. For more settings, please refer to the supplementary material.

self-optimize based on rendering loss during training, the given sampling positions become more precise, resulting in finer details than Monoangelo.

4.3 Ablations

The experiment was conducted on UrbanScene3D. The qualitative and quantitative evaluation results are shown in Fig 7 and Table 2. **LG Guider.** We freeze the parameters of the LG Guider (Freeze Prop) after step2 training, the sampling position given by LG Guider can no longer vary. When the parameters of the LG Guider are not affected by the rendering loss of step 3, we found that the ambiguity is somewhat alleviated. This because the LG Guider has already learned the geometric information at step1. However, LG Guider loses its ability to refine its sampling field during training, resulting in over-smoothing.

Non occupancy loss. L_{nocc} is designed to prevent the new surface from appearing in areas where σ should be smaller according to the reliable geometric information when we take rendering loss at optimization step 3. When L_{nocc} is removed, we can see that the scene has some raised surfaces at corner regions which is easily suffers the ambiguity.

Training strategy. We also carry out the experiments to prove that our training strategy is effective and outperforms other strategies mentioned in Fig 3. Please refer to our supplementary material.

Table 3: Quantitative results of the ablation study on the UrbanScene3D dataset.

Method	CD ↓	Acc ₉₅ ↓	Comp ₉₅ ↓	Overall ₉₅ ↓
SciArt				
Base	1.3938	0.3319	0.5813	0.4566
Freeze Prop	1.4585	0.3736	0.6982	0.5359
No L_{nocc}	1.1322	<u>0.2980</u>	0.4649	0.3815
Full	<u>1.0574</u>	0.2990	<u>0.4138</u>	<u>0.3564</u>
Polytech				
Base	1.1029	0.2989	0.3969	0.3479
Freeze Prop	0.8350	0.2218	0.3182	0.2700
No L_{nocc}	0.6782	<u>0.1749</u>	0.2120	0.1935
Full	<u>0.6593</u>	0.1763	<u>0.2086</u>	<u>0.1925</u>

5 Limitations and Future work

Due to the high reconstruction accuracy of MegaSurf, seams are usually imperceptible when assembling all the blocks together. However, in some cases, seams may still be perceptible. Applying mesh refinement to the assembled model can effectively overcome this issue. Recently, 3D Gaussian-based surface reconstruction [1, 2, 8, 9, 51] has become a research hotspot in surface reconstruction due to the fast training capability. While the reconstruction quality needs further improvement. MegaSurf’s training strategy can be adapted to Gaussian-based methods to enhance reconstruction accuracy. Additionally, Gaussian splatting requires a large number of 3D Gaussians to represent fine geometric surfaces, especially in large-scale scenes, which places high demands on GPU memory. Therefore, using implicit encoding for Gaussians or combining with SDFs may be promising future directions.

6 Conclusion

We introduced MegaSurf, a novel Learnable Sampling Guided surface reconstruction approach for reconstructing large-scale scenes. To accelerate the process, we developed a Fast PatchMatch MVS module that efficiently propagates SFM information to surrounding areas, yielding high-confidence geometric cues. Furthermore, we introduced the Learnable Geometric Guider (LG Guider) to learn a sampling field from these reliable geometric cues, which can be continuously refined through rendering loss minimization. To address the challenge of shape-radiance ambiguity, we employed a Divide-and-Conquer training strategy to harmonize the LG Guider and the radiance field, resulting in high-fidelity reconstructions. Extensive experiments on large-scale datasets demonstrate the superior performance of our method.

Acknowledgments

This work is partially supported by the National Natural Science Foundation of China (No. 62372336 and No. 61972298) and Wuhan University Huawei GeoInformatics Innovation Lab.

References

- [1] Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haoamin Liu, Hujun Bao, and Guofeng Zhang. 2024. PGSR: Planar-based Gaussian Splatting for Efficient and High-Fidelity Surface Reconstruction. *arXiv preprint arXiv:2406.06521* (2024).
- [2] Pinxuan Dai, Jiamin Xu, Wenxiang Xie, Xinguo Liu, Huamin Wang, and Weiwei Xu. 2024. High-quality surface reconstruction using gaussian surfels. In *ACM SIGGRAPH 2024 Conference Papers*. 1–11.
- [3] Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. 2022. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems* 35 (2022), 3403–3416.
- [4] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. 2015. Massively parallel multiview stereopsis by surface normal diffusion. In *Proceedings of the IEEE International Conference on Computer Vision*. 873–881.
- [5] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. 2020. 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*. 2495–2504.
- [6] Jianfei Guo, Nianchen Deng, Xinyang Li, Yeqi Bai, Botian Shi, Chiyu Wang, Chenjing Ding, Dongliang Wang, and Yikang Li. 2023. StreetSurf: Extending Multi-view Implicit Surface Reconstruction to Street Views. *arXiv preprint arXiv:2306.04988* (2023).
- [7] Heiko Hirschmüller. 2005. Accurate and efficient stereo processing by semi-global matching and mutual information. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, Vol. 2. IEEE, 807–814.
- [8] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2024. 2D Gaussian Splatting for Geometrically Accurate Radiance Fields. *arXiv preprint arXiv:2403.17888* (2024).
- [9] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* 42, 4 (2023), 1–14.
- [10] Jingliang Li, Zhengda Lu, Yiqun Wang, Ying Wang, and Jun Xiao. 2022. Dsmvnet: Unsupervised multi-view stereo via depth synthesis. In *Proceedings of the 30th ACM International Conference on Multimedia*. 5593–5601.
- [11] Zhuopeng Li, Lu Li, and Jianke Zhu. 2023. Read: Large-scale neural scene rendering for autonomous driving. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 1522–1529.
- [12] Zongcheng Li, Xiaoxiao Long, Yusen Wang, Tuo Cao, Wenping Wang, Fei Luo, and Chunxia Xiao. 2023. NeTO: neural reconstruction of transparent objects with self-occlusion aware refraction-tracing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 18547–18557.
- [13] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. 2023. Neuralangelo: High-Fidelity Neural Surface Reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8456–8465.
- [14] Jie Liao, Yanping Fu, Qingan Yan, Fei Luo, and Chunxia Xiao. 2021. Adaptive depth estimation for pyramid multi-view stereo. *Computers & Graphics* 97 (2021), 268–278.
- [15] Jie Liao, Yanping Fu, Qingan Yan, and Chunxia Xiao. 2019. Pyramid multi-view stereo with local consistency. In *Computer Graphics Forum*, Vol. 38. Wiley Online Library, 335–346.
- [16] Liqiang Lin, Yilin Liu, Yue Hu, Xingguang Yan, Ke Xie, and Hui Huang. 2022. Capturing, Reconstructing, and Simulating: the UrbanScene3D Dataset. In *ECCV*. 93–109.
- [17] Tianqi Liu, Xinyi Ye, Weiyue Zhao, Zhiyu Pan, Min Shi, and Zhiguo Cao. 2023. When epipolar constraint meets non-local operators in multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 18088–18097.
- [18] William E Lorensen and Harvey E Cline. 1987. Marching cubes: A high resolution 3D surface construction algorithm. *ACM siggraph computer graphics* 21, 4 (1987), 163–169.
- [19] Chunjie Luo, Fei Luo, Yusen Wang, Enxu Zhao, and Chunxia Xiao. 2024. DLCA-Recon: Dynamic Loose Clothing Avatar Reconstruction from Monocular Videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 3963–3971.
- [20] Zeyu Ma, Zachary Teed, and Jia Deng. 2022. Multiview stereo with cascaded epipolar raft. In *European Conference on Computer Vision*. Springer, 734–750.
- [21] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*. Springer, 405–421.
- [22] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–15.
- [23] Konstantinos Rematas, Andrew Liu, Pratul P. Srinivasan, Jonathan T. Barron, Andrea Tagliasacchi, Tom Funkhouser, and Vittorio Ferrari. 2022. Urban Radiance Fields. *CVPR* (2022).
- [24] Radu Alexandru Rosu and Sven Behnke. 2023. Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8466–8475.
- [25] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. 2016. Pixelwise view selection for unstructured multi-view stereo. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14*. Springer, 501–518.
- [26] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. 2006. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, Vol. 1. IEEE, 519–528.
- [27] Jiaming Sun, Xi Chen, Qianqian Wang, Zhengqi Li, Hadar Averbuch-Elor, Xiaowei Zhou, and Noah Snavely. 2022. Neural 3d reconstruction in the wild. In *ACM SIGGRAPH 2022 Conference Proceedings*. 1–9.
- [28] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. 2022. Blocknerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8248–8258.
- [29] Haimen Turki, Deva Ramanan, and Mahadev Satyanarayanan. 2022. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12922–12931.
- [30] Chen Wang, Xian Wu, Yuan-Chen Guo, Song-Hai Zhang, Yu-Wing Tai, and Shi-Min Hu. 2022. NeRF-SR: High-Quality Neural Radiance Fields using Super-sampling. In *Proceedings of the 30th ACM International Conference on Multimedia*. 6445–6454.
- [31] Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and Wenping Wang. 2022. Neuris: Neural reconstruction of indoor scenes using normal priors. In *European Conference on Computer Vision*. Springer, 139–155.
- [32] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689* (2021).
- [33] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. 2023. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3295–3306.
- [34] Yusen Wang, Zongcheng Li, Yu Jiang, Kaixuan Zhou, Tuo Cao, Yanping Fu, and Chunxia Xiao. 2022. NeuralRoom: Geometry-Constrained Neural Implicit Surfaces for Indoor Scene Reconstruction. *ACM Transactions on Graphics (TOG)* 41, 6 (2022), 1–15.
- [35] Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. 2022. Hf-neus: Improved surface reconstruction using high-frequency details. *Advances in Neural Information Processing Systems* 35 (2022), 1966–1978.
- [36] Zian Wang, Tianchang Shen, Jun Gao, Shengyu Huang, Jacob Munkberg, Jon Hasselgren, Zan Gojic, Wenzheng Chen, and Sanja Fidler. 2023. Neural Fields meet Explicit Geometric Representations for Inverse Rendering of Urban Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8370–8380.
- [37] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. 2021. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5610–5619.
- [38] Tong Wu, Jiaqi Wang, Xingang Pan, XU Xudong, Christian Theobalt, Ziwei Liu, and Dahua Lin. 2022. Voxurf: Voxel-based Efficient and Accurate Neural Surface Reconstruction. In *The Eleventh International Conference on Learning Representations*.
- [39] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. 2022. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *European conference on computer vision*. Springer, 106–122.
- [40] Luoyuan Xu, Tao Guan, Yuesong Wang, Yawei Luo, Zhuo Chen, Wenkai Liu, and Wei Yang. 2022. Self-supervised multi-view stereo via adjacent geometry guided volume completion. In *Proceedings of the 30th ACM International Conference on Multimedia*. 2202–2210.
- [41] Linning Xu, Yuanbo Xiangli, Sida Peng, Xingang Pan, Nanxuan Zhao, Christian Theobalt, Bo Dai, and Dahua Lin. 2023. Grid-guided Neural Radiance Fields for Large Urban Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8296–8306.
- [42] Qingshan Xu, Weihang Kong, Wenbing Tao, and Marc Pollefeys. 2022. Multi-scale geometric consistency guided and planar prior assisted multi-view stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 4 (2022), 4945–4963.
- [43] Qingshan Xu and Wenbing Tao. 2019. Multi-scale geometric consistency guided multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5483–5492.
- [44] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. 2022. Point-NeRF: Point-based Neural Radiance Fields.

- arXiv preprint arXiv:2201.08845* (2022).
- [45] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. 2018. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 767–783.
- [46] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems* 34 (2021), 4805–4815.
- [47] Lior Yariv, Peter Hedman, Christian Reiser, Dor Verbin, Pratul P. Srinivasan, Richard Szeliski, Jonathan T. Barron, and Ben Mildenhall. 2023. BakedSDF: Meshing Neural SDFs for Real-Time View Synthesis. *arXiv* (2023).
- [48] Zehao Yu, Anpei Chen, Bozidar Antic, Songyou Peng, Apratim Bhattacharyya, Michael Niemeyer, Siyu Tang, Torsten Sattler, and Andreas Geiger. 2022. SDFS-tudio: A Unified Framework for Surface Reconstruction. <https://github.com/autonomousvision/sdfstudio>
- [49] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. 2022. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems* 35 (2022), 25018–25032.
- [50] Junyi Zeng, Chong Bao, Rui Chen, Zilong Dong, Guofeng Zhang, Hujun Bao, and Zhaopeng Cui. 2023. Mirror-NeRF: Learning Neural Radiance Fields for Mirrors with Whitted-Style Ray Tracing. In *Proceedings of the 31st ACM International Conference on Multimedia*. 4606–4615.
- [51] Baowen Zhang, Chuan Fang, Rakesh Shrestha, Yixun Liang, Xiaoxiao Long, and Ping Tan. 2024. RaDe-GS: Rasterizing Depth in Gaussian Splatting. *arXiv preprint arXiv:2406.01467* (2024).
- [52] Jingyang Zhang, Yao Yao, Shiwei Li, Zixin Luo, and Tian Fang. 2020. Visibility-aware multi-view stereo network. *arXiv preprint arXiv:2008.07928* (2020).